

Adaptive VEM for variable data: convergence and optimality

L. Beirão da Veiga ^{*1}, C. Canuto ^{†2}, R. H. Nochetto ^{‡3}, G. Vacca ^{§4}, and M. Verani ^{¶5}

¹Dipartimento di Matematica e Applicazioni, Università degli Studi di Milano
Bicocca, Via Roberto Cozzi 55 - 20125 Milano, Italy

²Dipartimento di Scienze Matematiche G.L. Lagrange, Politecnico di Torino, Corso
Duca degli Abruzzi 24 - 10129 Torino, Italy

³Department of Mathematics and Institute for Physical Science and Technology,
University of Maryland, College Park - 20742, MD, USA

⁴Dipartimento di Matematica, Università degli Studi di Bari, Via Edoardo Orabona 4
- 70125 Bari, Italy

⁵MOX-Laboratory for Modeling and Scientific Computing, Dipartimento di
Matematica, Politecnico di Milano, Piazza Leonardo da Vinci 32 - 20133 Milano, Italy

February 28, 2023

Abstract

We design an adaptive virtual element method (AVEM) of lowest order over triangular meshes with hanging nodes in 2d, which are treated as polygons. AVEM hinges on the stabilization-free a posteriori error estimators recently derived in [8]. The crucial property, that also plays a central role in this paper, is that the stabilization term can be made arbitrarily small relative to the a posteriori error estimators upon increasing the stabilization parameter. Our AVEM concatenates two modules, **GALERKIN** and **DATA**. The former deals with piecewise constant data and is shown in [8] to be a contraction between consecutive iterates. The latter approximates general data by piecewise constants to a desired accuracy. AVEM is shown to be convergent and quasi-optimal, in terms of error decay versus degrees of freedom, for solutions and data belonging to appropriate approximation classes. Numerical experiments illustrate the interplay between these two modules and provide computational evidence of optimality.

1 Introduction

Virtual element methods (VEMs) are a new paradigm for the conforming discretization of partial differential equations (PDEs) over polytopal meshes. They were introduced a few

^{*}lourengo.beirao@unimib.it

[†]claudio.canuto@polito.it

[‡]rhn@math.umd.edu

[§]giuseppe.vacca@uniba.it

[¶]marco.verani@polimi.it

years ago and have seen a rapid development with an increasing number of applications ever since [5, 6, 7]. Virtual element functions are continuous piecewise polynomials on the skeleton of the polytopal mesh and are extended inside the elements in a convenient way that avoids their explicit manipulation. This flexibility allows for global regularity, say continuity in the context of second order PDEs, but requires dealing with projection operators and stabilization of the resulting discrete bilinear form to be coercive (or more generally to satisfy a discrete inf-sup condition). If the PDE has variable data $\mathcal{D} = (A, c, f)$, as in our prototype boundary value problem

$$-\nabla \cdot (A\nabla u) + cu = f \quad \text{in } \Omega, \quad u = 0 \quad \text{on } \partial\Omega, \quad (1.1)$$

then \mathcal{D} has to be further approximated to formulate the discrete counterpart of (1.1). This is well understood in the a priori analysis of VEMs, which deliver optimal convergence rates under minimal regularity assumptions on \mathcal{D} and for rather simple and practical choices of the stabilization term.

The a posteriori error analysis of VEMs approximations of (1.1) initiated in [3, 14, 9], along with suitable upper and lower error estimates for variable data \mathcal{D} . The stabilization term and the residual error estimator of [14], which is the one more relevant to us, turn out to be of the same order but the former is not bounded above by the energy error. This is problematic to study convergence of any adaptive VEM (AVEM for short). We have recently tackled this crucial issue in [8] and shown that the stabilization term can be made arbitrarily small relative to the error estimator upon increasing the stabilization parameter. This property is valid in 2d on newest-vertex bisection meshes made of triangles with hanging nodes and a fixed maximal *global index*, which limits the level of hanging nodes. Hence, triangles with multiple nodes are viewed as polygons for the VEM approach. This severe mesh restriction is crucial to relate the actual VEM mesh \mathcal{T} with the largest conforming submesh \mathcal{T}^0 of \mathcal{T} and their approximation properties. Moreover, this leads to stabilization-free a posteriori error estimates, derived in [8], and facilitates the convergence analysis of AVEM, which is the ultimate objective of this paper. We are not aware of similar studies for AVEM even though convergence is a fundamental mathematical question of practical significance.

In contrast, the convergence and optimality analyses of adaptive finite element methods (AFEMs) constitute a mature research field for elliptic PDEs such as (1.1); we refer to the surveys [18, 19] as well as [15] for details. A common approach in the AFEM literature is to assume that the linear and bilinear forms associated with (1.1) can be computed exactly. The role of quadrature is not assessed a posteriori and, as a consequence, the resulting AFEMs are not fully practical unless data \mathcal{D} is piecewise polynomial. This leads to the usual one-loop AFEMs which iterate the modules

$$\text{SOLVE} \longrightarrow \text{ESTIMATE} \longrightarrow \text{MARK} \longrightarrow \text{REFINE}. \quad (1.2)$$

A valid and practical alternative is to first approximate \mathcal{D} by piecewise polynomials to a desired accuracy, and next run (1.2) for such approximate data to achieve a comparable level of precision. This two-step AFEM was first proposed by R. Stevenson [20], and further explored in [11, 16].

Dealing with approximate data \mathcal{D} is inherent to the formulation of VEMs and their basic definition. It is thus natural in this context to think of two-step AVEMs. This is precisely our intent in this paper, in which we design an AVEM for (1.1) in two stages. We first assume that \mathcal{D} is piecewise constant and introduce a one-step AVEM, the so-called **GALERKIN**

module, which is shown in [8] to possess a contraction property between consecutive adaptive iterations. We next consider variable data \mathcal{D} and design a two-step AVEM that consists of a concatenation of the modules **DATA** and **GALERKIN** in the spirit of [11, 16, 20]. Given an initial mesh \mathcal{T}_0 and parameters $\varepsilon_0, \omega > 0$, AVEM sets $k = 0$ and iterates

$$\begin{aligned} [\widehat{\mathcal{T}}_k, \widehat{\mathcal{D}}_k] &= \text{DATA}(\mathcal{T}_k, \mathcal{D}, \omega \varepsilon_k) \\ [\mathcal{T}_{k+1}, u_{k+1}] &= \text{GALERKIN}(\widehat{\mathcal{T}}_k, \widehat{\mathcal{D}}_k, \varepsilon_k) \\ \varepsilon_{k+1} &= \frac{1}{2}\varepsilon_k; \quad k \leftarrow k + 1 \end{aligned}$$

The module **DATA** approximates $\mathcal{D} = (A, c, f)$ in the spaces $((L^\infty(\Omega))^{2 \times 2}, L^\infty(\Omega), L^2(\Omega))$ by piecewise constant data $\widehat{\mathcal{D}}_k$ on an admissible refinement $\widehat{\mathcal{T}}_k$ of \mathcal{T}_k to accuracy $\omega \varepsilon_k$. The pair $(\widehat{\mathcal{T}}_k, \widehat{\mathcal{D}}_k)$ is taken by **GALERKIN** to run an inner loop, with piecewise constant data $\widehat{\mathcal{D}}_k$ and initial mesh $\widehat{\mathcal{T}}_k$, that creates the next mesh-solution pair $(\mathcal{T}_{k+1}, u_{k+1})$. The module **GALERKIN** stops as soon as the error tolerance ε_k is reached, which takes a finite number of iterations because **GALERKIN** is a contraction between consecutive iterates. It is worth noticing that, in the absence of this stopping test, **GALERKIN** would converge to the solution of (1.1) corresponding to the perturbed data $\widehat{\mathcal{D}}_k$, which is not the desired solution u of (1.1). The relative resolution of the modules **DATA** and **GALERKIN** is critical and is governed by the parameter $\omega > 0$. In our numerical experiments we observe that $\omega = 1$ is an adequate choice.

It is clear from its definition that this two-step AVEM converges. Concerning its optimality, we show that the number of iterations of **GALERKIN** is independent of the iteration counter k and its complexity is dictated by the approximation classes of the solution u and data \mathcal{D} . This requires ω to be sufficiently small, or equivalently that the perturbed solution of (1.1) with data $\widehat{\mathcal{D}}_k$ is much closer to u than the error tolerance ε_k ; this is in the spirit of [11, 20]. We also prove that the complexity of **DATA** is given by suitable approximation classes of $\mathcal{D} = (A, c, f)$ in the spaces $((L^\infty(\Omega))^{2 \times 2}, L^\infty(\Omega), L^2(\Omega))$. Altogether, this yields the following optimal decay estimate for the energy error in terms of the number of degrees of freedom $\#\mathcal{T}_k$

$$|u - u_k|_{1,\Omega} \leq C(u, \mathcal{D}) (\#\mathcal{T}_k)^{-s}, \quad (1.3)$$

where $s > 0$ is the worse decay rate between those of the near-best approximations errors for u and for \mathcal{D} ; typically $s = \frac{1}{2}$ in dimension 2.

This paper is organized as follows. We present the weak formulation of (1.1) in Section 2 and recall the VEM basic ingredients in Section 3. We discuss VEM for piecewise constant data in Section 4, including the stabilization-free a posteriori error estimates from [8]. In Section 5 we design **GALERKIN**, and recall its fundamental contraction property from [8]. We deal with variable data in Section 6, which entails a perturbation estimate for (1.1), the design of **DATA**, and eventually of **AVEM** for general data. Section 7 analyzes the computational cost of **GALERKIN**, showing that the number of sub-iterations inside a call to **GALERKIN** is uniformly bounded. Section 8 is devoted to the study of the quasi-optimality of AVEM: approximation classes for the solution and data are introduced, and the rate decay of the error in the energy norm versus the number of degrees of freedom is estimated in terms of these classes. Section 9 completes the analysis, with the study of the decay of data approximation errors. We document the interplay between the modules **DATA** and **GALERKIN** with several illuminating numerical experiments in Section 10. It is important to realize that for mesh refinement to maintain bounded global indices, and thus admissible meshes, further refinement beyond

the marked elements might be necessary. In Section 11 we design and study a procedure to make meshes admissible in the sense that the global index is uniformly bounded for all k . This procedure hinges on the bisection algorithm and is of somewhat intrinsic interest. We prove that it is optimal in terms of degrees of freedom, very much in the spirit of the completion algorithm for conforming bisection meshes by Binev, Dahmen, and DeVore [10]; see also [18, 19, 21]. We finally draw conclusions in Section 12.

2 The continuous problem

In a polygonal domain $\Omega \subset \mathbb{R}^2$, consider the second-order Dirichlet boundary-value problem

$$-\nabla \cdot (A\nabla u) + cu = f \quad \text{in } \Omega, \quad u = 0 \quad \text{on } \partial\Omega, \quad (2.1)$$

with data $\mathcal{D} = (A, c, f)$, where $A \in (L^\infty(\Omega))^{2 \times 2}$ is symmetric and uniformly positive-definite in Ω , $c \in L^\infty(\Omega)$ is non-negative in Ω , and $f \in L^2(\Omega)$. The variational formulation of the problem is

$$u \in \mathbb{V} : \quad \mathcal{B}(u, v) = (f, v)_\Omega \quad \forall v \in \mathbb{V}, \quad (2.2)$$

with $\mathbb{V} := H_0^1(\Omega)$ and $\mathcal{B}(u, v) := a(u, v) + m(u, v)$ where

$$a(u, v) := \int_\Omega (A\nabla u) \cdot \nabla v, \quad m(u, v) := \int_\Omega c u v$$

are the bilinear forms associated with (2.1). Let $\|\cdot\| = \sqrt{\mathcal{B}(\cdot, \cdot)}$ be the energy norm, which satisfies

$$c_{\mathcal{B}} |v|_{1,\Omega}^2 \leq \|v\|^2 \leq c^{\mathcal{B}} |v|_{1,\Omega}^2 \quad \forall v \in \mathbb{V}, \quad (2.3)$$

for suitable constants $0 < c_{\mathcal{B}} \leq c^{\mathcal{B}}$.

3 VEM preliminaries

In view of the adaptive discretization of the problem, let us fix an initial conforming partition \mathcal{T}_0 of $\bar{\Omega}$ made of triangular elements. Let us denote by \mathcal{T} any refinement of \mathcal{T}_0 obtained by a finite number of successive *newest-vertex bisections* [10, 18, 19, 21]; the triangulation \mathcal{T} need not be conforming, since hanging nodes may be generated by the refinement. Let \mathcal{N} denote the set of nodes of \mathcal{T} , i.e., the collection of all vertices of the triangles in \mathcal{T} ; a node $z \in \mathcal{N}$ is *proper* if it is a vertex of all triangles containing it; otherwise, it is a *hanging node*. Thus, $\mathcal{N} = \mathcal{P} \cup \mathcal{H}$ is partitioned into the union of the set \mathcal{P} of proper nodes and the set \mathcal{H} of hanging nodes.

Given an element $E \in \mathcal{T}$, let \mathcal{N}_E be the set of nodes sitting on ∂E ; it contains the three vertices and, possibly, some hanging nodes. If the cardinality $|\mathcal{N}_E| = 3$, E is said a *proper triangle* of \mathcal{T} ; if $|\mathcal{N}_E| > 3$, then according to the VEM philosophy E is not viewed as a triangle, but as a polygon having $|\mathcal{N}_E|$ edges, some of which are placed consecutively on the same line; the set of all edges of E is denoted by \mathcal{E}_E . Note that if $e \subset \partial E \cap \partial E'$, then it is an edge for both elements; consequently, it is meaningful to define the *skeleton* of the triangulation \mathcal{T} by setting $\mathcal{E} = \mathcal{E}_{\mathcal{T}} := \bigcup_{E \in \mathcal{T}} \mathcal{E}_E$. Throughout the paper, we will set $h_E = |E|^{1/2}$ for an element and $h_e = |e|$ for an edge.

The concept of *global index* of a hanging node, introduced in [8], will be crucial in the sequel. To define it, let us first observe that any hanging node $\mathbf{x} \in \mathcal{H}$ has been obtained through a newest-vertex bisection by halving an edge of a triangle in the preceding triangulation; denoting by $\mathbf{x}', \mathbf{x}'' \in \mathcal{N}$ the endpoints of such edge, let us set $\mathbf{B}(\mathbf{x}) = \{\mathbf{x}', \mathbf{x}''\}$.

Definition 3.1 (Global index of a node and a partition). *The global index λ of a node $\mathbf{x} \in \mathcal{N}$ is recursively defined as follows:*

- If $\mathbf{x} \in \mathcal{P}$, then set $\lambda(\mathbf{x}) := 0$;
- If $\mathbf{x} \in \mathcal{H}$, with $\mathbf{x}', \mathbf{x}'' \in \mathbf{B}(\mathbf{x})$, then set $\lambda(\mathbf{x}) := \max(\lambda(\mathbf{x}'), \lambda(\mathbf{x}'')) + 1$.

The global index of the partition \mathcal{T} is defined as $\Lambda_{\mathcal{T}} := \max_{\mathbf{x} \in \mathcal{N}} \lambda(\mathbf{x})$.

Definition 3.2 (Λ -admissible partitions). *Given a constant $\Lambda \geq 1$, a non-conforming partition \mathcal{T} is said to be Λ -admissible if*

$$\Lambda_{\mathcal{T}} \leq \Lambda.$$

Starting from the initial conforming partition \mathcal{T}_0 (which is trivially Λ -admissible), all the subsequent non-conforming partitions generated by the module **REFINE** in the sequel will remain Λ -admissible due to the algorithm **CREATE_ADMISSIBLE_CHAIN** studied in Section 11. We refer to [12] for a similar algorithm in the context of dG approximations.

Remark 3.3. The condition that \mathcal{T} is Λ -admissible has the following implications for each element $E \in \mathcal{T}$:

- If $L \subset \partial E$ is one of the three sides of the triangle E , then L may contain at most $2^\Lambda - 1$ hanging nodes; consequently, $|\mathcal{N}_E| \leq 3 \cdot 2^\Lambda$.
- If $e \subset \partial E$ is any edge, then $h_e \simeq h_E$, where the hidden constants only depend on the shape of the initial triangulation \mathcal{T}_0 and possibly on Λ .

In the following C will denote a generic positive constant independent of the mesh \mathcal{T} but which may depend on Ω , on the initial partition \mathcal{T}_0 , on the data \mathcal{D} and on the constant Λ (cf. Definition 3.2) and that may change at each occurrence, whereas the symbol \lesssim will denote a bound up to C .

3.1 VEM spaces and projectors

Although the results of the present paper apply to a wider set of VEM spaces [5, 1, 7], we prefer to focus the attention on the so-called *enhanced* VEM space. We will be brief and refer to [8] for a more detailed description which adopts the same notation. We start with the projector $\Pi_E^\nabla : H^1(E) \rightarrow \mathbb{P}_1(E)$, which is defined by the conditions

$$(\nabla(v - \Pi_E^\nabla v), \nabla q_1)_E = 0 \quad \forall q_1 \in \mathbb{P}_1(E), \quad \int_{\partial E} (v - \Pi_E^\nabla v) = 0. \quad (3.1)$$

To introduce the space of discrete functions in Ω associated with \mathcal{T} , for each element $E \in \mathcal{T}$ we define

$$\begin{aligned} \mathbb{V}_{\partial E} &:= \{v \in \mathcal{C}^0(\partial E) : v|_e \in \mathbb{P}_1(e) \quad \forall e \in \mathcal{E}_E\}, \\ \mathbb{V}_E &:= \{v \in H^1(E) : v|_{\partial E} \in \mathbb{V}_{\partial E}, \Delta v \in \mathbb{P}_1(E), \int_E (v - \Pi_E^\nabla v) q_1 = 0 \quad \forall q_1 \in \mathbb{P}_1(E)\}. \end{aligned} \quad (3.2)$$

Obviously $\mathbb{P}_1(E) \subseteq \mathbb{V}_E$ and, if E is a proper triangle, then $\mathbb{V}_E = \mathbb{P}_1(E)$. Once the local spaces \mathbb{V}_E are defined, we introduce the global discrete space

$$\mathbb{V}_{\mathcal{T}} := \{v \in \mathbb{V} : v|_E \in \mathbb{V}_E \quad \forall E \in \mathcal{T}\}. \quad (3.3)$$

Note that functions in $\mathbb{V}_{\mathcal{T}}$ are piecewise affine on the skeleton \mathcal{E} and are globally continuous. A *set of degrees of freedom* for the space $\mathbb{V}_{\mathcal{T}}$ is given by the pointwise evaluation at all (internal) mesh vertices.

We also define the subspace of continuous, piecewise affine functions on \mathcal{T}

$$\mathbb{V}_{\mathcal{T}}^0 := \{v \in \mathbb{V} : v|_E \in \mathbb{P}_1(E) \quad \forall E \in \mathcal{T}\} \subseteq \mathbb{V}_{\mathcal{T}}. \quad (3.4)$$

This subspace was crucial in [8] to get a stabilization-free a posteriori error estimate, and will play an essential role in this paper as well to remove the stabilization term from several estimates.

The discretization of Problem (2.1) will involve the following global projection operators

$$\Pi_{\mathcal{T}}^{\nabla} : \mathbb{V}_{\mathcal{T}} \rightarrow \mathbb{P}_1(\mathcal{T}), \quad \mathcal{I}_{\mathcal{T}} : \mathbb{V}_{\mathcal{T}} \rightarrow \mathbb{P}_1(\mathcal{T}), \quad \Pi_{\mathcal{T}}^0 : L^2(\Omega) \rightarrow \mathbb{P}_1(\mathcal{T}), \quad (3.5)$$

where $\mathbb{P}_1(\mathcal{T})$ denotes the space of (discontinuous) piecewise linear polynomials over \mathcal{T} . We define these operators in terms of their local counterparts. In fact, for each element $E \in \mathcal{T}$, $\Pi_{\mathcal{T}}^{\nabla}$ restricts to the local elliptic projection operator Π_E^{∇} in (3.1), $\mathcal{I}_{\mathcal{T}}$ restricts to the local Lagrange interpolation operator $\mathcal{I}_E : \mathbb{V}_E \rightarrow \mathbb{P}_1(E)$ at the vertices of E , and Π_E^0 restricts to the local L^2 -orthogonal projection operator $\Pi_E^0 : L^2(E) \rightarrow \mathbb{P}_1(E)$. It turns out that $\Pi_E^0 = \Pi_E^{\nabla}$ on \mathbb{V}_E , because of the definition (3.2) of the space \mathbb{V}_E , and that Π_E^{∇} is computable on \mathbb{V}_E in terms of the degrees of freedom [5, 8]. Furthermore, in view of the definition (3.4) of $\mathbb{V}_{\mathcal{T}}^0$, $\Pi_{\mathcal{T}}^{\nabla} v = \mathcal{I}_{\mathcal{T}} v = v$ for all $v \in \mathbb{V}_{\mathcal{T}}^0$.

4 A Virtual Element Method with piecewise constant data

In this section we briefly summarize the definition and certain properties of the virtual element discretization of (2.2) introduced in [8] under the following assumption.

Assumption 4.1 (coefficients and right-hand side of the equation). *The data $\mathcal{D} = (A, c, f)$ in (2.1) are constant in each element E of \mathcal{T} .*

For any $E \in \mathcal{T}$ we use the following notation: $A_E = A|_E \in \mathbb{R}^{2 \times 2}$, $c_E = c|_E \in \mathbb{R}$, $f_E = f|_E \in \mathbb{R}$.

4.1 The discrete problem

Under the above assumption, we define $a_{\mathcal{T}}, m_{\mathcal{T}} : \mathbb{V}_{\mathcal{T}} \times \mathbb{V}_{\mathcal{T}} \rightarrow \mathbb{R}$ by

$$\begin{aligned} a_{\mathcal{T}}(v, w) &:= \sum_{E \in \mathcal{T}} a_E(v, w), & a_E(v, w) &:= \int_E (A_E \nabla \Pi_E^{\nabla} v) \cdot \nabla \Pi_E^{\nabla} w, \\ m_{\mathcal{T}}(v, w) &:= \sum_{E \in \mathcal{T}} m_E(v, w), & m_E(v, w) &:= c_E \int_E \Pi_E^{\nabla} v \Pi_E^{\nabla} w. \end{aligned} \quad (4.1)$$

Next, for any $E \in \mathcal{T}$, we introduce the stabilization symmetric bilinear form $s_E : \mathbb{V}_E \times \mathbb{V}_E \rightarrow \mathbb{R}$

$$s_E(v, w) = \sum_{i=1}^{\mathcal{N}_E} v(\mathbf{x}_i) w(\mathbf{x}_i), \quad (4.2)$$

with $\{\mathbf{x}_i\}_{i=1}^{\mathcal{N}_E}$ denoting the nodes of E . This form controls the kernel of a_E on \mathbb{V}_E/\mathbb{R} because it satisfies

$$c_s |v|_{1,E}^2 \leq s_E(v, v) \leq C_s |v|_{1,E}^2 \quad \forall v \in \mathbb{V}_E/\mathbb{R}, \quad (4.3)$$

for constants $C_s \geq c_s > 0$ independent of E ; for a proof of (4.3) we refer to [2, 13]. Other choices for the stabilization form are available in the literature [2, 13] and the results presented here easily extend to such cases. With the local form s_E at hand, we define the local and global stabilization forms

$$\begin{aligned} S_E(v, w) &:= s_E(v - \mathcal{I}_E v, w - \mathcal{I}_E w) \quad \forall v, w \in \mathbb{V}_E, \\ S_{\mathcal{T}}(v, w) &:= \sum_{E \in \mathcal{T}} S_E(v, w) \quad \forall v, w \in \mathbb{V}_{\mathcal{T}}. \end{aligned} \quad (4.4)$$

Note that from (4.3) we obtain

$$S_{\mathcal{T}}(v, v) \simeq |v - \mathcal{I}_{\mathcal{T}} v|_{1,\mathcal{T}}^2 \quad \forall v \in \mathbb{V}_{\mathcal{T}}, \quad (4.5)$$

where $|\cdot|_{1,\mathcal{T}}$ denotes the broken H^1 -seminorm over the mesh \mathcal{T} .

Finally, for all $v, w \in \mathbb{V}_{\mathcal{T}}$ we define the complete bilinear form

$$\mathcal{B}_{\mathcal{T}} : \mathbb{V}_{\mathcal{T}} \times \mathbb{V}_{\mathcal{T}} \rightarrow \mathbb{R}, \quad \mathcal{B}_{\mathcal{T}}(v, w) := a_{\mathcal{T}}(v, w) + m_{\mathcal{T}}(v, w) + \gamma S_{\mathcal{T}}(v, w), \quad (4.6)$$

where $\gamma \geq \gamma_0$ for some fixed $\gamma_0 > 0$ is a stabilization constant independent of \mathcal{T} . The following properties are an easy consequence of the definitions and bounds outlined above.

Lemma 4.2 (properties of bilinear forms). *The following properties are valid*

- For any $v \in \mathbb{V}_{\mathcal{T}}$ and any $w \in \mathbb{V}_{\mathcal{T}}^0$, it holds

$$a_{\mathcal{T}}(v, w) = a(v, w), \quad m_{\mathcal{T}}(v, w) = m(v, w), \quad S_{\mathcal{T}}(v, w) = 0. \quad (4.7)$$

- The form $\mathcal{B}_{\mathcal{T}}$ satisfies

$$b |v|_{1,\Omega}^2 \leq \mathcal{B}_{\mathcal{T}}(v, v), \quad |\mathcal{B}_{\mathcal{T}}(v, w)| \leq B |v|_{1,\Omega} |w|_{1,\Omega}, \quad \forall v, w \in \mathbb{V}_{\mathcal{T}}, \quad (4.8)$$

with continuity and coercivity constants $B \geq b > 0$ independent of the triangulation \mathcal{T} .

Recalling (4.6), direct consequence of (4.7) is the following consistency result:

$$\mathcal{B}_{\mathcal{T}}(v, w) = B(v, w) \quad \forall v \in \mathbb{V}_{\mathcal{T}}, \forall w \in \mathbb{V}_{\mathcal{T}}^0. \quad (4.9)$$

We now have all the ingredients to set the Galerkin discretization of Problem (2.1): find

$$u_{\mathcal{T}} \in \mathbb{V}_{\mathcal{T}} : \quad \mathcal{B}_{\mathcal{T}}(u_{\mathcal{T}}, v) = \mathcal{F}_{\mathcal{T}}(v) \quad \forall v \in \mathbb{V}_{\mathcal{T}}, \quad (4.10)$$

with discrete loading term

$$\mathcal{F}_{\mathcal{T}}(v) := \sum_{E \in \mathcal{T}} f_E \int_E \Pi_E^{\nabla} v \quad \forall v \in H_0^1(\Omega). \quad (4.11)$$

Combining (4.8) with the Lax-Milgram Lemma, we obtain existence, uniqueness and stability of the solution $u_{\mathcal{T}}$ of (4.10). Moreover, $u_{\mathcal{T}}$ satisfies the following orthogonality condition in the subspace $\mathbb{V}_{\mathcal{T}}^0$ [5, 8].

Lemma 4.3 (Galerkin quasi-orthogonality). *The solutions u of (2.2) and $u_{\mathcal{T}}$ of (4.10) satisfy*

$$\mathcal{B}(u - u_{\mathcal{T}}, v) = 0 \quad \forall v \in \mathbb{V}_{\mathcal{T}}^0. \quad (4.12)$$

4.2 An a posteriori error estimator

Since we are interested in building adaptive discretizations, we rely on a posteriori error control. Hereafter we present the residual-type a posteriori estimator introduced in [8] as a variant of the one in [14]. To this end, recalling that $\mathcal{D} = (A, c, f)$ denotes the set of piecewise constant data, for any $v \in \mathbb{V}_{\mathcal{T}}$ and any element E let us define the internal residual over E

$$r_{\mathcal{T}}(E; v, \mathcal{D}) := f_E - c_E \Pi_E^{\nabla} v. \quad (4.13)$$

Similarly, for any two elements $E_1, E_2 \in \mathcal{T}$ sharing an edge $e \in \mathcal{E}_{E_1} \cap \mathcal{E}_{E_2}$, let us define the jump residual over e

$$j_{\mathcal{T}}(e; v, \mathcal{D}) := \llbracket A \nabla \Pi_{\mathcal{T}}^{\nabla} v \rrbracket_e = (A_{E_1} \nabla \Pi_{E_1}^{\nabla} v|_{E_1}) \cdot \mathbf{n}_1 + (A_{E_2} \nabla \Pi_{E_2}^{\nabla} v|_{E_2}) \cdot \mathbf{n}_2, \quad (4.14)$$

where \mathbf{n}_i denotes the unit normal vector to e pointing outward with respect to E_i ; set $j_{\mathcal{T}}(e; v, \mathcal{D}) = 0$ if $e \subset \partial\Omega$. Then, taking into account Remark 3.3, we define the local residual estimator associated with E

$$\eta_{\mathcal{T}}^2(E; v, \mathcal{D}) := h_E^2 \|r_{\mathcal{T}}(E; v, \mathcal{D})\|_{0,E}^2 + \frac{1}{2} \sum_{e \in \mathcal{E}_E} h_e \|j_{\mathcal{T}}(e; v, \mathcal{D})\|_{0,e}^2. \quad (4.15)$$

The residual estimator localized on some subset $\mathcal{S} \subseteq \mathcal{T}$ is

$$\eta_{\mathcal{T}}^2(\mathcal{S}; v, \mathcal{D}) := \sum_{E \in \mathcal{S}} \eta_{\mathcal{T}}^2(E; v, \mathcal{D}) \quad (4.16)$$

and the global residual estimator is

$$\eta_{\mathcal{T}}^2(v, \mathcal{D}) := \eta_{\mathcal{T}}^2(\mathcal{T}; v, \mathcal{D}) = \sum_{E \in \mathcal{T}} \eta_{\mathcal{T}}^2(E; v, \mathcal{D}). \quad (4.17)$$

Upper and lower a posteriori bounds of the energy error are provided by the following result, whose proof can be found in [8, Proposition 4.1 and Corollary 4.3].

Proposition 4.4 (a posteriori error estimates). *There exist constants $C_{apost} > c_{apost} > 0$ depending on Λ and \mathcal{D} but independent of u , \mathcal{T} , $u_{\mathcal{T}}$ and the stabilization parameter γ , such that*

$$\begin{aligned} |u - u_{\mathcal{T}}|_{1,\Omega}^2 &\leq C_{apost} (\eta_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}) + S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}})) , \\ c_{apost} \eta_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}) &\leq |u - u_{\mathcal{T}}|_{1,\Omega}^2 + S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}}) . \end{aligned} \quad (4.18)$$

The stabilization term $S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}})$ and residual estimator $\eta_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D})$ are, unfortunately, of the same order [8, Section 4.1]. However, such difficulty is handled by the following crucial result, proved in [8, Proposition 4.4], which relies on the subspace $\mathbb{V}_{\mathcal{T}}^0$ and Lemma 4.3 (Galerkin quasi-orthogonality). This shows the importance of $\mathbb{V}_{\mathcal{T}}^0$.

Proposition 4.5 (bound of the stabilization term by the residual). *There exists a constant $C_B > 0$, depending on Λ but independent of \mathcal{T} , $u_{\mathcal{T}}$ and the stabilization parameter γ , such that*

$$\gamma^2 S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}}) \leq C_B \eta_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}). \quad (4.19)$$

Combining (4.18) and (4.19) gives rise to the following fundamental estimate [8, Corollary 4.5].

Theorem 4.6 (stabilization-free a posteriori error estimates). *Assume that the stabilization parameter γ is chosen to satisfy $\gamma^2 > \frac{C_B}{c_{\text{apost}}}$. Then it holds*

$$C_L \eta_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}) \leq |u - u_{\mathcal{T}}|_{1,\Omega}^2 \leq C_U \eta_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}), \quad (4.20)$$

with $C_L = c_{\text{apost}} - C_B \gamma^{-2}$ and $C_U = C_{\text{apost}} (1 + C_B \gamma^{-2})$.

5 AVEM for piecewise constant data

In this section, we recall from [8] the Adaptive Virtual Element Method (AVEM) for approximating (2.2) under Assumption 4.1, together with its convergence property. In particular, AVEM for piecewise constant data is realized by a call to the module **GALERKIN** described hereafter. Given a Λ -admissible input mesh $\widehat{\mathcal{T}}$, piecewise constant input data \mathcal{D} on $\widehat{\mathcal{T}}$ and a tolerance $\varepsilon > 0$, the module

$$[\mathcal{T}, u_{\mathcal{T}}] = \text{GALERKIN}(\widehat{\mathcal{T}}, \mathcal{D}, \varepsilon) \quad (5.1)$$

produces a Λ -admissible bisection refinement \mathcal{T} of $\widehat{\mathcal{T}}$ and the Galerkin approximation $u_{\mathcal{T}} \in \mathbb{V}_{\mathcal{T}}$ to the solution u of problem (2.1) with piecewise constant data \mathcal{D} , such that

$$\|u - u_{\mathcal{T}}\| \leq C_G \varepsilon, \quad (5.2)$$

with $C_G = \sqrt{c^{\mathcal{B}} C_U}$, where $c^{\mathcal{B}}$ is defined in (2.3) and C_U is defined in (4.20). This is obtained by iterating the classical paradigm

$$\text{SOLVE} \longrightarrow \text{ESTIMATE} \longrightarrow \text{MARK} \longrightarrow \text{REFINE} \quad (5.3)$$

producing a sequence of Λ -admissible meshes $\{\mathcal{T}_k\}_{k \geq 0}$, with $\mathcal{T}_0 = \widehat{\mathcal{T}}$, and associated Galerkin solutions $u_k \in \mathbb{V}_{\mathcal{T}_k}$ to the problem (2.1) with data \mathcal{D} . The iteration stops as soon as $\eta_{\mathcal{T}_k}(u_k, \mathcal{D}) \leq \varepsilon$, which is possible thanks to the convergence result stated in Theorem 5.2 below.

The modules in (5.3) are defined as follows: given piecewise constant data \mathcal{D} on \mathcal{T}_0 ,

- $[u_{\mathcal{T}}] = \text{SOLVE}(\mathcal{T}, \mathcal{D})$ produces the Galerkin solution on the mesh \mathcal{T} for data \mathcal{D} ;

- $[\{\eta_{\mathcal{T}}(\cdot; u_{\mathcal{T}}, \mathcal{D})\}] = \text{ESTIMATE}(\mathcal{T}, u_{\mathcal{T}})$ computes the local residual estimators (4.15) on the mesh \mathcal{T} , which depend on the Galerkin solution $u_{\mathcal{T}}$ and data \mathcal{D} ;
- $[\mathcal{M}] = \text{MARK}(\mathcal{T}, \{\eta_{\mathcal{T}}(\cdot; u_{\mathcal{T}}, \mathcal{D})\}, \theta)$ implements the Dörfler criterion [17], precisely for a given parameter $\theta \in (0, 1)$ an almost minimal set $\mathcal{M} \subset \mathcal{T}$ is found such that

$$\theta \eta_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}) \leq \eta_{\mathcal{T}}^2(\mathcal{M}; u_{\mathcal{T}}, \mathcal{D}); \quad (5.4)$$

- $[\mathcal{T}_*] = \text{REFINE}(\mathcal{T}, \mathcal{M})$ produces a Λ -admissible refinement \mathcal{T}_* of \mathcal{T} , obtained by newest-vertex bisection of all the elements in \mathcal{M} and, possibly, some other elements.

In the procedure **REFINE**, non-admissible hanging nodes, i.e., hanging nodes with global index larger than Λ , might be created while refining elements in \mathcal{M} through newest-vertex bisection. Thus, in order to obtain a Λ -admissible partition \mathcal{T}_* , **REFINE** possibly refines other elements in \mathcal{T} . This is accomplished by applying to each $E \in \mathcal{M}$ a procedure, termed **CREATE_ADMISSIBLE_CHAIN**(\mathcal{T}, E), which identifies and refines a chain of elements starting at E , thereby creating a Λ -admissible partition. The loop is as follows:

```

 $[\mathcal{T}_*] = \text{REFINE}(\mathcal{T}, \mathcal{M})$ 
  for  $E \in \mathcal{M} \cap \mathcal{T}$ 
     $[\mathcal{T}] = \text{CREATE\_ADMISSIBLE\_CHAIN}(\mathcal{T}, E)$ 
  end for
  return( $\mathcal{T}$ )

```

Due to the technical nature of the procedure **CREATE_ADMISSIBLE_CHAIN**, we postpone its description and analysis to Section 11.1. We state now a complexity estimate for **REFINE**, whose proof is given at the end of that section. This result is fundamental for our optimality analysis of AVEM in Section 8 and is similar in spirit to the original estimate for the bisection method by Binev, Dahmen, and DeVore [10]; see also [18, 19, 21].

Theorem 5.1 (complexity of **REFINE**). *Let \mathcal{T}_0 be an initial mesh with suitable initial labeling. Let \mathcal{T}_k be a Λ -admissible refinement of \mathcal{T}_0 by newest-vertex bisection created by successive calls $\mathcal{T}_{j+1} = \text{REFINE}(\mathcal{T}_j, \mathcal{M}_j)$ for $0 \leq j \leq k-1$. Then there exists a universal constant $C_0 > 0$, solely depending on \mathcal{T}_0 and its labeling, such that*

$$\#\mathcal{T}_k - \#\mathcal{T}_0 \leq C_0 \sum_{j=0}^{k-1} \#\mathcal{M}_j. \quad (5.5)$$

We point out that a different procedure, termed **MAKE_ADMISSIBLE**, was used in [8] to generate a Λ -admissible refinement. While the implementation of this procedure is simpler than the one in **CREATE_ADMISSIBLE_CHAIN**, and works well in practice, only the latter guarantees the validity of the bound (5.5).

At last, we state the following convergence result for **GALERKIN** (cf. [8, Theorem 5.1]) with piecewise constant data.

Theorem 5.2 (convergence of **GALERKIN**). *There exist constants $\beta > 0$ and $\alpha \in (0, 1)$ such that, choosing the stabilization parameter $\gamma > 0$ sufficiently large in the Definition 4.6, the approximations $u_k \in \mathbb{V}_{\mathcal{T}_k}$ defined in **GALERKIN** satisfy*

$$\|u - u_k\|^2 + \beta \eta_{\mathcal{T}_k}^2(u_k, \mathcal{D}) \lesssim \alpha^k, \quad k \geq 0. \quad (5.6)$$

6 AVEM for general data

In this section we describe the two-step AVEM for general (non-piecewise constant) data and discuss its convergence properties. We first state the regularity of data.

Assumption 6.1 (regularity of data). *The data $\mathcal{D} = (A, c, f)$ satisfies*

$$\mathcal{D} \in C^0(\mathcal{T}_0; \mathbb{R}^{2 \times 2}) \times L^\infty(\Omega) \times L^2(\Omega),$$

where $C^0(\mathcal{T}_0; \mathbb{R}^{2 \times 2})$ denotes the space of piecewise uniformly continuous tensor fields over \mathcal{T}_0 .

We will see below that the regularity of c and f can be weakened, but not that of A unless we proceed as in [11]. We could assume $c \in L^q(\Omega)$ for $1 < q < \infty$ and $f \in H^{-1}(\Omega)$, but we will not pursue this regularity much further. We begin with a perturbation result for the solution of the exact problem.

6.1 Data perturbation

Let $\widehat{\mathcal{D}} = (\widehat{A}, \widehat{c}, \widehat{f})$ be the element-by-element average of $\mathcal{D} = (A, c, f)$ over a partition \mathcal{T} of Ω , namely

$$\widehat{A}|_E := A_E = \frac{1}{|E|} \int_E A \quad \widehat{c}|_E := c_E = \frac{1}{|E|} \int_E c \quad \widehat{f}|_E := f_E = \frac{1}{|E|} \int_E f \quad \forall E \in \mathcal{T}. \quad (6.1)$$

If $\alpha > 0$ is the smallest eigenvalue of $A = A(\mathbf{x})$ for all $\mathbf{x} \in \Omega$, then for any $\xi \in \mathbb{R}^2$

$$\xi \cdot A(\mathbf{x}) \xi \geq \alpha |\xi|^2 \quad \forall \mathbf{x} \in \Omega \quad \Rightarrow \quad \xi \cdot A_E \xi \geq \alpha |\xi|^2 \quad \forall E \in \mathcal{T},$$

whence the smallest eigenvalue $\widehat{\alpha}$ of \widehat{A} satisfies $\widehat{\alpha} \geq \alpha$; thus \widehat{A} is uniformly SPD in Ω . We view $\widehat{\mathcal{D}}$ as a perturbation of \mathcal{D} and consider the corresponding bilinear form $\widehat{\mathcal{B}}(\cdot, \cdot) = \widehat{a}(\cdot, \cdot) + \widehat{m}(\cdot, \cdot)$, with

$$\widehat{a}(u, v) = \int_\Omega \widehat{A} \nabla u \cdot \nabla v, \quad \widehat{m}(u, v) = \int_\Omega \widehat{c} u v \quad \forall u, v \in \mathbb{V},$$

and perturbed problem

$$\widehat{u} \in \mathbb{V} : \quad \widehat{\mathcal{B}}(\widehat{u}, v) = (\widehat{f}, v) \quad \forall v \in \mathbb{V}. \quad (6.2)$$

Lemma 6.2 (continuous dependence on data). *There exists a constant $C > 0$, depending on Ω and the mesh shape-regularity, such that for any $1 < q \leq \infty$ and $s \in [0, 1]$ satisfying $s < 2(q-1)/q$ it holds*

$$|u - \widehat{u}|_{1, \Omega} \leq \frac{1}{\alpha} |u|_{1, \Omega} \left(\|A - \widehat{A}\|_{L^\infty(\Omega)} + \frac{Cq}{2q-2-sq} \|h^s(c - \widehat{c})\|_{L^q(\Omega)} \right) + C \|h(f - \widehat{f})\|_{L^2(\Omega)}, \quad (6.3)$$

where the mesh density h of \mathcal{T} is the piecewise constant function satisfying $h|_E = h_E$ for all $E \in \mathcal{T}$.

Proof. We write the difference between (2.2) and (6.2) as follows:

$$\int_{\Omega} \left[(\hat{A} \nabla(u - \hat{u})) \cdot \nabla v + \hat{c}(u - \hat{u})v \right] = \int_{\Omega} \left[((\hat{A} - A) \nabla u) \cdot \nabla v + (\hat{c} - c)uv \right] + \int_{\Omega} (f - \hat{f})v.$$

Since (\hat{c}, \hat{f}) are L^2 -projections of (c, f) on piecewise constants over \mathcal{T} , we readily obtain

$$\int_{\Omega} \left[(\hat{A} \nabla(u - \hat{u})) \cdot \nabla v + \hat{c}(u - \hat{u})v \right] = \int_{\Omega} \left[((\hat{A} - A) \nabla u) \cdot \nabla v + (\hat{c} - c)(uv - \overline{uv}) \right] + \int_{\Omega} (f - \hat{f})(v - \overline{v}),$$

where the overbars denote the piecewise constant averages over \mathcal{T} . Taking the test function $v = u - \hat{u} \in \mathbb{V}$, and using the relation $\hat{\alpha} \geq \alpha > 0$ for the smallest eigenvalues of \hat{A} and A , standard arguments yield

$$\begin{aligned} \alpha \|\nabla v\|_{L^2(\Omega)}^2 &\leq \|(\hat{A} - A) \nabla u\|_{L^2(\Omega)} \|\nabla v\|_{L^2(\Omega)} \\ &\quad + \sum_{E \in \mathcal{T}} h_E^s \|c - \hat{c}\|_{L^q(E)} |uv|_{W_{q'}^s(E)} + \sum_{E \in \mathcal{T}} h_E \|f - \hat{f}\|_{L^2(E)} \|\nabla v\|_{L^2(E)}, \end{aligned} \quad (6.4)$$

where we set $q' = q/(q - 1)$. We now focus on the more involved mass term. We start by observing that, by a standard Hölder inequality on sequences and Sobolev embeddings,

$$\sum_{E \in \mathcal{T}} h_E^s \|c - \hat{c}\|_{L^q(E)} |uv|_{W_{q'}^s(E)} \leq \|\mathbf{h}^s(c - \hat{c})\|_{L^q(\Omega)} |uv|_{W_{q'}^s(\Omega)} \leq C \|\mathbf{h}^s(c - \hat{c})\|_{L^q(\Omega)} |uv|_{W_{p'}^1(\Omega)}, \quad (6.5)$$

where $1/p' = 1/q' - s/2$ and C depends on Ω . Consequently, for r satisfying $1/r + 1/2 = 1/p'$, we get

$$|uv|_{W_{p'}^1(\Omega)} = |u \nabla v + v \nabla u|_{L^{p'}(\Omega)} \leq \|u\|_{L^r(\Omega)} \|v\|_{H^1(\Omega)} + \|v\|_{L^r(\Omega)} \|u\|_{H^1(\Omega)}.$$

Combining the definitions of r and p' we easily obtain the explicit expression $r = 2q'/(2 - sq') = 2q/(2q - 2 - sq)$. Since $r \in [1, \infty)$, the Sobolev embedding $H^1(\Omega) \subseteq L^r(\Omega)$ and previous bound yield

$$|uv|_{W_{p'}^1(\Omega)} \leq Cr \|\nabla u\|_{L^2(\Omega)} \|\nabla v\|_{L^2(\Omega)} \quad (6.6)$$

with C depending on Ω . Inequalities (6.4), (6.5), (6.6) give

$$\begin{aligned} \alpha \|\nabla v\|_{L^2(\Omega)}^2 &\leq \|(\hat{A} - A) \nabla u\|_{L^2(\Omega)} \|\nabla v\|_{L^2(\Omega)} \\ &\quad + Cr \|\mathbf{h}^s(c - \hat{c})\|_{L^q(\Omega)} \|\nabla u\|_{L^2(\Omega)} \|\nabla v\|_{L^2(\Omega)} + C \|\mathbf{h}(f - \hat{f})\|_{L^2(\Omega)} \|\nabla v\|_{L^2(\Omega)}, \end{aligned}$$

from which we immediately deduce the asserted estimate (6.3). \square

Remark 6.3. The bound $|u|_{1,\Omega} \leq \frac{1}{\alpha} \|f\|_{H^{-1}(\Omega)}$ allows us to rewrite (6.3) in terms of data

$$|u - \hat{u}|_{1,\Omega} \leq \frac{1}{\alpha^2} \|f\|_{H^{-1}(\Omega)} \left(\|A - \hat{A}\|_{L^\infty(\Omega)} + \frac{Cq}{2q - 2 - sq} \|\mathbf{h}^s(c - \hat{c})\|_{L^q(\Omega)} \right) + C \|\mathbf{h}(f - \hat{f})\|_{L^2(\Omega)}. \quad (6.7)$$

Remark 6.4. A pair of relevant choices for q in Lemma 6.2 are $q = \infty$, which allows us to take $s = 1$, and $q = 2$, which allows us to take any value of s strictly smaller than one. Values of $q \leq 2$ can be taken, but note that the largest possible exponent s tends to zero as $q \rightarrow 1$.

6.2 The module DATA: piecewise constant approximation of data

Given $\mathcal{D} = (A, c, f)$ satisfying Assumption 6.1, a mesh \mathcal{T} and a tolerance ε , the module

$$[\widehat{\mathcal{T}}, \widehat{\mathcal{D}}] = \text{DATA}(\mathcal{T}, \mathcal{D}, \varepsilon) \quad (6.8)$$

produces a Λ -admissible bisection refinement $\widehat{\mathcal{T}}$ of \mathcal{T} and a piecewise constant approximation $\widehat{\mathcal{D}} = (\widehat{A}, \widehat{c}, \widehat{f})$ of \mathcal{D} over $\widehat{\mathcal{T}}$ such that

$$\|A - \widehat{A}\|_{L^\infty(\Omega)} + \|\mathbf{h}(c - \widehat{c})\|_{L^\infty(\Omega)} + \|\mathbf{h}(f - \widehat{f})\|_{L^2(\Omega)} \leq \varepsilon, \quad (6.9)$$

which controls the perturbation error according to Lemma 6.2 (continuous dependence on data). In view of (6.9), for any $E \in \widehat{\mathcal{T}}$ we introduce the following local data error estimators

$$\zeta_{\widehat{\mathcal{T}}}(E; A) := \|A - \widehat{A}\|_{L^\infty(E)}, \quad \zeta_{\widehat{\mathcal{T}}}(E; c) := h_E \|c - \widehat{c}\|_{L^\infty(E)}, \quad \zeta_{\widehat{\mathcal{T}}}(E; f) := h_E \|f - \widehat{f}\|_{L^2(E)}, \quad (6.10)$$

and the global data error estimators

$$\zeta_{\widehat{\mathcal{T}}}(A) := \|A - \widehat{A}\|_{L^\infty(\Omega)}, \quad \zeta_{\widehat{\mathcal{T}}}(c) := \|\mathbf{h}(c - \widehat{c})\|_{L^\infty(\Omega)}, \quad \zeta_{\widehat{\mathcal{T}}}(f) := \|\mathbf{h}(f - \widehat{f})\|_{L^2(\Omega)}, \quad (6.11)$$

and

$$\zeta_{\widehat{\mathcal{T}}}(\mathcal{D}) := \zeta_{\widehat{\mathcal{T}}}(A) + \zeta_{\widehat{\mathcal{T}}}(c) + \zeta_{\widehat{\mathcal{T}}}(f). \quad (6.12)$$

The data error reduction is obtained by iterating the following loop

$$\text{PROJECT} \longrightarrow \text{ESTIMATE_DATA} \longrightarrow \text{MARK_DATA} \longrightarrow \text{REFINE}, \quad (6.13)$$

which produces a sequence of Λ -admissible meshes $\{\widehat{\mathcal{T}}_j\}_{j \geq 0}$, with $\widehat{\mathcal{T}}_0 = \mathcal{T}$, and associated piecewise constant data $\widehat{\mathcal{D}}_j = (\widehat{A}_j, \widehat{c}_j, \widehat{f}_j)$ w.r.t. $\widehat{\mathcal{T}}_j$, that approximates the exact data \mathcal{D} until a $k \geq 0$ is found that satisfies $\zeta_{\widehat{\mathcal{T}}_k}(\mathcal{D}) \leq \varepsilon$.

The modules in (6.13) are defined as follows:

- $[\widehat{\mathcal{D}}] = \text{PROJECT}(\mathcal{T}, \mathcal{D})$ computes the element-by-element average $\widehat{\mathcal{D}} = (\widehat{A}, \widehat{c}, \widehat{f})$ of \mathcal{D} over \mathcal{T} ;
- $\{\{\zeta_{\mathcal{T}}(\cdot; A)\}, \{\zeta_{\mathcal{T}}(\cdot; c)\}, \{\zeta_{\mathcal{T}}(\cdot; f)\}\} = \text{ESTIMATE_DATA}(\mathcal{T}, \mathcal{D}, \widehat{\mathcal{D}})$ computes the local data error estimators (6.10) on the mesh \mathcal{T} ;
- $[\mathcal{M}_{\mathcal{D}}] = \text{MARK_DATA}(\mathcal{T}, \{\zeta_{\mathcal{T}}(\cdot; A)\}, \{\zeta_{\mathcal{T}}(\cdot; c)\}, \{\zeta_{\mathcal{T}}(\cdot; f)\}, \theta, \varepsilon)$ implements the following marking criteria. For the diffusion and the reaction terms A and c we apply the *greedy* strategy that selects

$$\mathcal{M}_A := \{E \in \mathcal{T} : \zeta_{\mathcal{T}}(E; A) \geq \tfrac{1}{3}\varepsilon\}, \quad \mathcal{M}_c := \{E \in \mathcal{T} : \zeta_{\mathcal{T}}(E; c) \geq \tfrac{1}{3}\varepsilon\}.$$

For the load term f , which accumulates in ℓ^2 rather than ℓ^∞ , we first check if $\zeta_{\mathcal{T}}(f) \geq \tfrac{1}{3}\varepsilon$, and if so we apply a *pseudo-greedy* strategy that, given a parameter $\theta \in (0, 1)$, selects

$$\mathcal{M}_f := \{E \in \mathcal{T} : \zeta_{\mathcal{T}}(E; f) \geq \theta \max_{E' \in \mathcal{T}} \zeta_{\mathcal{T}}(E'; f)\}. \quad (6.14)$$

Finally, we let the marked set be $\mathcal{M}_{\mathcal{D}} := \mathcal{M}_A \cup \mathcal{M}_c \cup \mathcal{M}_f$. In Sect. 9, the optimality properties of the greedy and pseudo-greedy strategies will be assessed.

- $[\widehat{\mathcal{T}}] = \text{REFINE}(\mathcal{T}, \mathcal{M}_{\mathcal{D}})$ produces a Λ -admissible refinement $\widehat{\mathcal{T}}$ of \mathcal{T} , obtained by newest-vertex bisection of all the elements in $\mathcal{M}_{\mathcal{D}}$ and, possibly, some other elements. This is the same procedure described in Section 5, applied with \mathcal{M} replaced by $\mathcal{M}_{\mathcal{D}}$.

Altogether, if \widehat{u} denotes the exact solution of the perturbed problem (6.2) with the output data $\widehat{\mathcal{D}}$ from (6.8), in view of (2.3) and (6.7) with $q = \infty$ there exists a constant C_D depending on Ω , data \mathcal{D} , and the shape-regularity constant of \mathcal{T}_0 such that **DATA** delivers the error estimate

$$\|u - \widehat{u}\| \leq C_D \varepsilon. \quad (6.15)$$

6.3 Realization of AVEM

Hereafter, we propose an adaptive VEM (or AVEM) that concatenates the modules **DATA** and **GALERKIN** introduced in (6.8) and (5.1), respectively. Concerning the latter module, its input now is a mesh $\widehat{\mathcal{T}}$ and piecewise constant data $\widehat{\mathcal{D}}$ on $\widehat{\mathcal{T}}$, while its output is a bisection refinement \mathcal{T} of $\widehat{\mathcal{T}}$ and the corresponding Galerkin approximation $u_{\mathcal{T}}$ to the exact solution \widehat{u} of problem (2.1) with piecewise constant data $\widehat{\mathcal{D}}$. They satisfy (5.2), namely

$$\|\widehat{u} - u_{\mathcal{T}}\| \leq C_G \varepsilon. \quad (6.16)$$

The module AVEM. Given an initial tolerance $\varepsilon_0 > 0$, a target tolerance **tol** and initial mesh \mathcal{T}_0 , as well as a safety parameter $\omega \in (0, 1]$, AVEM consists of the two-step algorithm:

```

 $[\mathcal{T}, u_{\mathcal{T}}] = \text{AVEM}(\mathcal{T}_0, \varepsilon_0, \omega, \text{tol})$ 
 $k = 0$ 
while  $\varepsilon_k > \frac{1}{2} \text{tol}$ 
   $[\widehat{\mathcal{T}}_k, \widehat{\mathcal{D}}_k] = \text{DATA}(\mathcal{T}_k, \mathcal{D}, \omega \varepsilon_k)$ 
   $[\mathcal{T}_{k+1}, u_{k+1}] = \text{GALERKIN}(\widehat{\mathcal{T}}_k, \widehat{\mathcal{D}}_k, \varepsilon_k)$ 
   $\varepsilon_{k+1} = \frac{1}{2} \varepsilon_k$ 
   $k \leftarrow k + 1$ 
end while
return  $(\mathcal{T}_k, u_k)$ 

```

Proposition 6.5 (convergence of AVEM). *For each $k \geq 0$ the modules **DATA** and **GALERKIN** converge in a finite number of iterations. Moreover, there exists a constant C_* depending solely on \mathcal{T}_0 such that the output of $[\mathcal{T}_{k+1}, u_{k+1}] = \text{GALERKIN}(\widehat{\mathcal{T}}_k, \widehat{\mathcal{D}}_k, \varepsilon_k)$ satisfies $\|u - u_{k+1}\| \leq C_* \varepsilon_k$ for all $k \geq 0$. Therefore, AVEM stops after K iterations, and delivers the estimate*

$$\|u - u_{K+1}\| \leq C_* \text{tol}.$$

Proof. We recall that Assumption 6.1 guarantees that A is uniformly continuous in each element of the initial mesh \mathcal{T}_0 . Consequently, $\|A - \widehat{A}\|_{L^\infty(E)}$ can be made arbitrarily small upon reducing h_E for all $E \in \widehat{\mathcal{T}}_k$. Moreover, since $c \in L^\infty(\Omega)$ and $f \in L^2(\Omega)$ in view of

Assumption 6.1, the errors $\|\mathbf{h}(c - \widehat{c})\|_{L^\infty(\Omega)}$ and $\|\mathbf{h}(f - \widehat{f})\|_{L^2(\Omega)}$ can also be made arbitrarily small because of the factor \mathbf{h} . This implies that **DATA** converges to tolerance $\omega\varepsilon_k$ for every $k \geq 0$ in a finite number of steps. The same is valid for **GALERKIN**, this time due to Theorem 5.2 (convergence of **GALERKIN**), whence we deduce that each loop of **AVEM** requires finite iterations. Thus, the output u_{k+1} satisfies

$$\|u - u_{k+1}\| \leq \|u - \widehat{u}_k\| + \|\widehat{u}_k - u_{k+1}\| \leq (C_D + C_G)\varepsilon_k \quad \forall k \geq 0,$$

according to (6.15) with $\omega\varepsilon_k \leq \varepsilon_k$ and (6.16). Finally, **AVEM** terminates after K loops, where K satisfies $\frac{1}{2}\mathbf{tol} < \varepsilon_K \leq \mathbf{tol}$, and the asserted estimate holds with $C_* = C_D + C_G$. \square

This elementary proof gives neither information about the dependence of the number of sub-iterations within each loop of **AVEM** upon the iteration counter k , nor insight whether the error decays optimally in terms of degrees of freedom. Answers to these two questions will be provided in Section 7 and Sections 8 and 9, respectively.

7 Computational cost of GALERKIN

In the sequel, we aim at investigating the complexity of **GALERKIN** within the **AVEM** loops. To this end, we need some preparatory results. In order to facilitate the reader, we shall use the notation

- **exact.sol**(\cdot), to indicate the exact solution to the boundary-value problem (2.1) with data prescribed by the argument,
- **galerkin.sol**(\cdot, \cdot), to indicate the solution to the Galerkin problem (4.10) on the partition prescribed by the first argument, with data prescribed by the second argument.

Furthermore, for any $k \in \mathbb{N}$, let $(\widehat{\mathcal{T}}_k, \widehat{\mathcal{D}}_k)$ and $(\mathcal{T}_{k+1}, u_{k+1})$, respectively, be the outputs of the module **DATA** and module **GALERKIN** at iteration k of **AVEM**. Then referring to (3.3), (3.5), (4.17), (6.1), we set the following notations:

mesh	VEM space	projection	estimator	piecewise constant data
$\widehat{\mathcal{T}}_k$	$\widehat{\mathbb{V}}_k := \mathbb{V}_{\widehat{\mathcal{T}}_k}$	$\widehat{\Pi}_k^\nabla := \Pi_{\widehat{\mathcal{T}}_k}^\nabla$	$\widehat{\eta}_k := \eta_{\widehat{\mathcal{T}}_k}$	$\widehat{\mathcal{D}}_k := (\widehat{A}_k, \widehat{c}_k, \widehat{f}_k)$
\mathcal{T}_k	$\mathbb{V}_k := \mathbb{V}_{\mathcal{T}_k}$	$\Pi_k^\nabla := \Pi_{\mathcal{T}_k}^\nabla$	$\eta_k := \eta_{\mathcal{T}_k}$	$\widehat{\mathcal{D}}_{k-1} := (\widehat{A}_{k-1}, \widehat{c}_{k-1}, \widehat{f}_{k-1})$.

Lemma 7.1 (uniform boundedness of u_k). *For any $k \geq 1$, let $u_k = \mathbf{galerkin.sol}(\mathcal{T}_k, \widehat{\mathcal{D}}_{k-1})$ be the output of the module **GALERKIN** at iteration $k - 1$. Then it holds*

$$|u_k|_{1,\Omega} \leq c_0 \|f\|_{0,\Omega} \tag{7.1}$$

for a constant $c_0 > 0$ independent of k .

Proof. Choosing $v = u_k = u_{\mathcal{T}_k}$ in (4.10) and noting that $\|\widehat{f}_{k-1}\|_{0,\Omega} \leq \|f\|_{0,\Omega}$, we get

$$\mathcal{B}_{\mathcal{T}_k}(u_k, u_k) = \mathcal{F}_{\mathcal{T}_k}(u_k) \leq \|f\|_{0,\Omega} \|\Pi_k^\nabla u_k\|_{0,\Omega}.$$

The result follows from the uniform H^1 -coercivity of the form $\mathcal{B}_{\mathcal{T}_k}$ and the H^1 -stability of the Π_k^∇ operator. \square

Lemma 7.2 (data perturbation of the error estimators). *For any $k \geq 1$, let (\mathcal{T}_k, u_k) be the output of the module GALERKIN at iteration $k-1$ of AVEM, i.e. $u_k = \text{galerkin.sol}(\mathcal{T}_k, \widehat{\mathcal{D}}_{k-1})$. Let $(\widehat{\mathcal{T}}_k, \widehat{\mathcal{D}}_k)$ be the output of the module DATA at iteration k of AVEM, and $u_{k,0} = \text{galerkin.sol}(\widehat{\mathcal{T}}_k, \widehat{\mathcal{D}}_k)$. Then it holds*

$$\widehat{\eta}_k^2(u_{k,0}, \widehat{\mathcal{D}}_k) \leq c_1 \eta_k^2(u_k, \widehat{\mathcal{D}}_{k-1}) + c_2 \epsilon_k^2 + c_3 |u_{k,0} - u_k|_{1,\Omega}^2 \quad (7.2)$$

for suitable positive constants c_1, c_2, c_3 .

Proof. We introduce the following notation

$$\begin{aligned} \widehat{r}_k &:= \widehat{f}_k - \widehat{c}_k \widehat{\Pi}_k^\nabla u_k & \widehat{j}_k &:= \llbracket \widehat{A}_k \nabla \widehat{\Pi}_k^\nabla u_k \rrbracket_{\widehat{\mathcal{E}}_k}, \\ r_k &:= \widehat{f}_{k-1} - \widehat{c}_{k-1} \Pi_k^\nabla u_k & j_k &:= \llbracket \widehat{A}_{k-1} \nabla \Pi_k^\nabla u_k \rrbracket_{\widehat{\mathcal{E}}_k}, \end{aligned}$$

and we observe that it holds

$$\widehat{r}_k = r_k + \widehat{c}_k (\Pi_k^\nabla u_k - \widehat{\Pi}_k^\nabla u_k) + (\widehat{f}_k - \widehat{f}_{k-1}) + (\widehat{c}_{k-1} - \widehat{c}_k) \Pi_k^\nabla u_k, \quad (7.3)$$

$$\widehat{j}_k = j_k + \llbracket \widehat{A}_k \nabla (\widehat{\Pi}_k^\nabla - \Pi_k^\nabla) u_k \rrbracket_{\widehat{\mathcal{E}}_k} + \llbracket (\widehat{A}_k - \widehat{A}_{k-1}) \nabla \Pi_k^\nabla u_k \rrbracket_{\widehat{\mathcal{E}}_k}. \quad (7.4)$$

We distinguish between refined and unrefined elements. Let us start from refined elements and let E be an element of \mathcal{T}_k which is split into $E_1, \dots, E_n \in \widehat{\mathcal{T}}_k$, where n depends on E and it holds $\min_{1 \leq i \leq n} h_{E_i} \leq h_E/2$. Hence, we have

$$\begin{aligned} \sum_{i=1}^n h_{E_i}^2 \|\widehat{r}_k\|_{E_i}^2 &\lesssim \sum_{i=1}^n h_{E_i}^2 \left(\|r_k\|_{E_i}^2 + \|\widehat{c}_k (\Pi_k^\nabla u_k - \widehat{\Pi}_k^\nabla u_k)\|_{E_i}^2 \right) \\ &\quad + \sum_{i=1}^n h_{E_i}^2 \left(\|\widehat{f}_k - \widehat{f}_{k-1}\|_{E_i}^2 + \|(\widehat{c}_{k-1} - \widehat{c}_k) \Pi_k^\nabla u_k\|_{E_i}^2 \right) =: I + II, \\ \sum_{i=1}^n \sum_{e \in \mathcal{E}_{E_i}} h_{E_i} \|\widehat{j}_k\|_e^2 &\lesssim \sum_{i=1}^n \sum_{e \in \mathcal{E}_{E_i}} \left(h_{E_i} \|j_k\|_e^2 + h_{E_i} \|\llbracket \widehat{A}_k \nabla (\widehat{\Pi}_k^\nabla - \Pi_k^\nabla) u_k \rrbracket_e\|_e^2 \right) \\ &\quad + \sum_{i=1}^n \sum_{e \in \mathcal{E}_{E_i}} \left(h_{E_i} \|\llbracket (\widehat{A}_k - \widehat{A}_{k-1}) \nabla \Pi_k^\nabla u_k \rrbracket_e\|_e^2 \right) =: III + IV. \end{aligned}$$

Adapting to I and III the same reasoning as in the proof of [8, Lemma 5.2] we get

$$\widehat{\eta}_k^2(E; u_k, \widehat{\mathcal{D}}_k) \lesssim \eta_k^2(E; u_k, \widehat{\mathcal{D}}_{k-1}) + S_{\mathcal{T}_k(E)}(u_k, u_k) + II + IV. \quad (7.5)$$

By employing [8, Lemma 5.3] we get

$$\widehat{\eta}_k^2(E; u_{k,0}, \widehat{\mathcal{D}}_k) \lesssim \eta_k^2(E; u_k, \widehat{\mathcal{D}}_{k-1}) + S_{\mathcal{T}_k(E)}(u_k, u_k) + |u_{k,0} - u_k|_{1,\mathcal{T}(E)}^2 + II + IV. \quad (7.6)$$

The sum $II + IV$ can be bounded using Hölder's inequality, the trace inequality together with (6.9), the stability property of Π^∇ and Lemma 7.1, obtaining

$$II + IV \lesssim \epsilon_k^2.$$

On unrefined elements E , we note that $\Pi_k^{\nabla,E} u_k = \widehat{\Pi}_k^{\nabla,E} u_k$. Hence, employing (7.3)-(7.4) together with [8, Lemma 5.3], and estimating the terms $II + IV$ as before, we have

$$\widehat{\eta}_k^2(E; u_{k,0}, \widehat{\mathcal{D}}_k) \lesssim \eta_k^2(E; u_k, \widehat{\mathcal{D}}_{k-1}) + \epsilon_k^2 + |u_{k,0} - u_k|_{1,E}^2. \quad (7.7)$$

Finally, summing over E and employing (4.19), we have

$$\hat{\eta}_k^2(u_{k,0}, \hat{\mathcal{D}}_k) \lesssim \eta_k^2(u_k, \hat{\mathcal{D}}_{k-1}) + \varepsilon_k^2 + |u_{k,0} - u_k|_{1,\Omega}^2. \quad (7.8)$$

□

Proposition 7.3 (computational cost of GALERKIN). *For any $k \in \mathbb{N}$, the number J_k of sub-iterations inside the call to GALERKIN at iteration k of AVEM is bounded independently of k .*

Proof. We proceed in several steps. For any $k \in \mathbb{N}$, let $(\hat{\mathcal{T}}_k, \hat{\mathcal{D}}_k)$ and $(\mathcal{T}_{k+1}, u_{k+1})$ be the output respectively of the module DATA and module GALERKIN at iteration k of AVEM. We will use the following functions:

$$\begin{aligned} \hat{u}_{k-1} &= \text{exact.sol}(\hat{\mathcal{D}}_{k-1}) \in \mathbb{V} & \hat{u}_k &= \text{exact.sol}(\hat{\mathcal{D}}_k) \in \mathbb{V} \\ u_k &= \text{galerkin.sol}(\mathcal{T}_k, \hat{\mathcal{D}}_{k-1}) \in \mathbb{V}_k & u_{k+1} &= \text{galerkin.sol}(\mathcal{T}_{k+1}, \hat{\mathcal{D}}_k) \in \mathbb{V}_{k+1} \\ u_k^{\text{en}} &= \text{galerkin.sol}(\hat{\mathcal{T}}_k, \hat{\mathcal{D}}_{k-1}) \in \hat{\mathbb{V}}_k & u_{k,0} &= \text{galerkin.sol}(\hat{\mathcal{T}}_k, \hat{\mathcal{D}}_k) \in \hat{\mathbb{V}}_k, \end{aligned} \quad (7.9)$$

where the suffix “en” stands for “enhanced” (i.e., u_k^{en} is computed with the same data as u_k , but on a finer mesh).

Step 1. Estimate of $|\hat{u}_k - u_{k,0}|_{1,\Omega}$. This is a consequence of the a posteriori error upper bound

$$|\hat{u}_k - u_{k,0}|_{1,\Omega}^2 \leq C_U \hat{\eta}_k^2(u_{k,0}, \hat{\mathcal{D}}_k)$$

given in Theorem 4.6.

Step 2. Estimate of $\hat{\eta}_k^2(u_{k,0}, \hat{\mathcal{D}}_k)$. Lemma 7.2 gives

$$\hat{\eta}_k^2(u_{k,0}, \hat{\mathcal{D}}_k) \leq c_1 \eta_k^2(u_k, \hat{\mathcal{D}}_{k-1}) + c_2 \varepsilon_k^2 + c_3 |u_{k,0} - u_k|_{1,\Omega}^2$$

which, in view of the input tolerance ε_k appearing in the module GALERKIN, implies

$$\hat{\eta}_k^2(u_{k,0}, \hat{\mathcal{D}}_k) \leq c_4 \varepsilon_k^2 + c_3 |u_{k,0} - u_k|_{1,\Omega}^2 \quad (7.10)$$

for some $c_4 > 0$. It remains to estimate $|u_{k,0} - u_k|_{1,\Omega}^2$ which, invoking the triangle inequality, reduces to

$$|u_{k,0} - u_k|_{1,\Omega} \lesssim |u_{k,0} - u_k^{\text{en}}|_{1,\Omega} + |u_k^{\text{en}} - u_k|_{1,\Omega}. \quad (7.11)$$

We observe that the difference $u_{k,0} - u_k^{\text{en}}$ between the two Galerkin solutions in $\hat{\mathbb{V}}_k$ is the solution of the following variational problem: for any $v \in \hat{\mathbb{V}}_k$ it holds

$$\begin{aligned} & \int_{\Omega} \hat{A}_{k-1} \nabla \hat{\Pi}_k^{\nabla} (u_{k,0} - u_k^{\text{en}}) \cdot \nabla \hat{\Pi}_k^{\nabla} v + \int_{\Omega} \hat{c}_{k-1} \hat{\Pi}_k^{\nabla} (u_{k,0} - u_k^{\text{en}}) \hat{\Pi}_k^{\nabla} v + S_{\hat{\mathcal{T}}_k} (u_{k,0} - u_k^{\text{en}}, v) = \\ & = \int_{\Omega} (\hat{A}_{k-1} - \hat{A}_k) \nabla \hat{\Pi}_k^{\nabla} u_{k,0} \cdot \nabla \hat{\Pi}_k^{\nabla} v + \int_{\Omega} (\hat{c}_{k-1} - \hat{c}_k) \hat{\Pi}_k^{\nabla} u_{k,0} \hat{\Pi}_k^{\nabla} v + \int_{\Omega} (\hat{f}_k - \hat{f}_{k-1}) \hat{\Pi}_k^{\nabla} v. \end{aligned}$$

Taking $v = u_{k,0} - u_k^{\text{en}}$, employing on the left-hand side the uniform coercivity of the discrete bilinear term, and using on the right-hand side the triangle inequality, the Cauchy-Schwarz inequality together with (6.9), and Lemma 7.1, we get

$$|u_{k,0} - u_k^{\text{en}}|_{1,\Omega} \leq c_5 (\varepsilon_k + \varepsilon_{k-1}) = 3c_5 \varepsilon_k \quad (7.12)$$

for a proper choice of $c_5 > 0$. In order to estimate $|u_k^{\text{en}} - u_k|_{1,\Omega}$, we preliminary note that $\widehat{\mathcal{T}}_k$ is a refinement of \mathcal{T}_k . Hence, invoking [8, Corollary 5.8] we have

$$\|\widehat{u}_{k-1} - u_k^{\text{en}}\|^2 + \|u_k^{\text{en}} - u_k\|^2 \leq (1 + 4\delta) \|\widehat{u}_{k-1} - u_k\|^2$$

which, in view of (2.3), yields

$$|u_k^{\text{en}} - u_k|_{1,\Omega} \leq c_6 |\widehat{u}_{k-1} - u_k|_{1,\Omega} \quad (7.13)$$

for some $c_6 > 0$. On the other hand, from Theorem 4.6, we have

$$|\widehat{u}_{k-1} - u_k|_{1,\Omega} \leq \sqrt{C_U} \eta_k(u_k, \widehat{\mathcal{D}}_{k-1}) \leq \sqrt{C_U} \varepsilon_{k-1} = 2\sqrt{C_U} \varepsilon_k. \quad (7.14)$$

Thus, from eqs. (7.11)-(7.14), we obtain

$$|u_{k,0} - u_k|_{1,\Omega} \leq (3c_5 + 2\sqrt{C_U} c_6) \varepsilon_k \quad (7.15)$$

and, employing (7.10), we arrive at

$$\widehat{\eta}_k^2(u_{k,0}, \widehat{\mathcal{D}}_k) \leq c_7 \varepsilon_k^2$$

for some $c_7 > 0$.

Step 3. Estimate of the total error $\xi_{\widehat{\mathcal{T}}_k}^2(u_{k,0})$, where, referring to Theorem 5.2, for any refinement \mathcal{T}_* of $\widehat{\mathcal{T}}_k$ and for any $v \in \mathbb{V}_{\mathcal{T}_*}$ we set

$$\xi_{\mathcal{T}_*}^2(v) := |\widehat{u}_k - v|_{1,\Omega}^2 + \beta \eta_{\mathcal{T}_*}^2(v, \widehat{\mathcal{D}}_k).$$

Because of Steps 1 and 2 we have

$$\xi_{\widehat{\mathcal{T}}_k}^2(u_{k,0}) \leq c_7 (C_U + \beta) \varepsilon_k^2 =: c_8 \varepsilon_k^2.$$

Step 4. Bound on J_k . Each consecutive iterate $(\mathcal{T}_{k,j}, u_{k,j})$ inside **GALERKIN** starting with $(\mathcal{T}_{k,0}, u_{k,0}) = (\widehat{\mathcal{T}}_k, u_{k,0})$ satisfies the contraction property in Theorem 5.2 (cf. [8, Theorem 5.1]). Therefore

$$\xi_{\mathcal{T}_{k,j}}^2(u_{k,j}) \lesssim \alpha^j \xi_{\widehat{\mathcal{T}}_k}^2(u_{k,0}) \leq \alpha^j c_9 \varepsilon_k^2,$$

for some $c_9 > 0$. Since J_k is the smallest value for which

$$\eta_{\mathcal{T}_{k,J_k}}(u_{k,J_k}, \widehat{\mathcal{D}}_k) \leq \varepsilon_k$$

we have

$$\eta_{\mathcal{T}_{k,J_k-1}}(u_{k,J_k-1}, \widehat{\mathcal{D}}_k) > \varepsilon_k.$$

Concatenating the last two ingredients gives

$$\varepsilon_k^2 \leq \frac{1}{\beta} \xi_{\mathcal{T}_{k,J_k-1}}^2(u_{k,J_k-1}) \leq \alpha^{J_k-1} \frac{c_9}{\beta} \varepsilon_k^2.$$

This in turn implies

$$\left(\frac{1}{\alpha}\right)^{J_k-1} \leq \frac{c_9}{\beta} \Rightarrow J_k \leq 1 + \frac{\log(c_9/\beta)}{\log(1/\alpha)} =: J.$$

We see that the upper bound J of J_k is independent of k . This concludes the proof. \square

8 Quasi-optimal cardinality of AVEM

The main purpose of this section is to prove, under suitable assumptions on the solution u and data \mathcal{D} , the bound (1.3) announced in the Introduction, namely the existence of constants $C(u, \mathcal{D}) > 0$ and $s \in (0, \frac{1}{2}]$ such that

$$|u - u_k|_{1,\Omega} \leq C(u, \mathcal{D}) (\#\mathcal{T}_k)^{-s}. \quad (8.1)$$

To this end, we introduce in Sect. 8.1 certain approximation classes for functions in \mathbb{V} and for data, tailored on the decomposition of Ω into Λ -admissible non-conforming partitions, and we assume that the solution and the data of Problem (2.1) belong to some of these classes. In Sect. 8.2, we investigate the approximability properties of certain perturbations of the exact solution, namely exact solutions of (2.1) with perturbed coefficients. Next, in Sect. 8.3, we consider a refinement \mathcal{T}_* of a partition \mathcal{T} , and give conditions under which an optimal Dörfler marking property holds. This allows us to prove in Sect. 8.4 an optimal estimate of the cardinality of the marked set in a call to GALERKIN. At last, in Sect. 8.5, we apply these results to establish the desired estimate on the rate of decay of the error produced by AVEM.

8.1 Approximation classes

We first introduce two families of approximation classes for a function $v \in \mathbb{V}$, and we show they coincide. Subsequently, we define approximation classes for the operator coefficients $A \in (L^\infty(\Omega))^{2 \times 2}$ and $c \in L^\infty(\Omega)$, and for the forcing $f \in L^2(\Omega)$.

8.1.1 Approximation classes for $v \in \mathbb{V}$

We start by defining the following quantity for $v \in \mathbb{V}$ and $v_{\mathcal{T}} \in \mathbb{V}_{\mathcal{T}}$

$$\mathbb{E}_{\mathcal{T}}^2(v, v_{\mathcal{T}}) := \|v - v_{\mathcal{T}}\|^2 + |v_{\mathcal{T}} - \mathcal{I}_{\mathcal{T}} v_{\mathcal{T}}|_{1,\mathcal{T}}^2. \quad (8.2)$$

It is worthy to observe that for $v_{\mathcal{T}}^0 \in \mathbb{V}_{\mathcal{T}}^0$ it obviously holds

$$\mathbb{E}_{\mathcal{T}}^2(v, v_{\mathcal{T}}^0) = \|v - v_{\mathcal{T}}^0\|^2. \quad (8.3)$$

Lemma 8.1 (quasi-best approximation). *Let u and $u_{\mathcal{T}}$ be the solutions of problem (2.2) and problem (4.10), respectively, with piecewise constant data. There exists a constant $C^\dagger > 0$, independent of u and the mesh \mathcal{T} , such that*

$$\mathbb{E}_{\mathcal{T}}^2(u, u_{\mathcal{T}}) \leq C^\dagger \mathbb{E}_{\mathcal{T}}^2(u, v_{\mathcal{T}}) \quad \forall v_{\mathcal{T}} \in \mathbb{V}_{\mathcal{T}}. \quad (8.4)$$

Proof. Let $\varepsilon_{\mathcal{T}} = u_{\mathcal{T}} - v_{\mathcal{T}}$. By the triangle inequality

$$\mathbb{E}_{\mathcal{T}}^2(u, u_{\mathcal{T}}) \leq 2 (\|u - v_{\mathcal{T}}\|^2 + |v_{\mathcal{T}} - \mathcal{I}_{\mathcal{T}} v_{\mathcal{T}}|_{1,\mathcal{T}}^2 + \|\varepsilon_{\mathcal{T}}\|^2 + |\varepsilon_{\mathcal{T}} - \mathcal{I}_{\mathcal{T}} \varepsilon_{\mathcal{T}}|_{1,\mathcal{T}}^2), \quad (8.5)$$

so that we only need to bound the last two terms. First by the coercivity of the discrete bilinear form, then by recalling the discrete (4.10) and continuous (2.2) weak problems, we obtain

$$\|\varepsilon_{\mathcal{T}}\|^2 + |\varepsilon_{\mathcal{T}} - \mathcal{I}_{\mathcal{T}} \varepsilon_{\mathcal{T}}|_{1,\mathcal{T}}^2 \leq C \mathcal{B}_{\mathcal{T}}(\varepsilon_{\mathcal{T}}, \varepsilon_{\mathcal{T}}) = C(\mathcal{F}_{\mathcal{T}}(\varepsilon_{\mathcal{T}}) - \mathcal{B}_{\mathcal{T}}(v_{\mathcal{T}}, \varepsilon_{\mathcal{T}})) = C(\mathcal{B}(u, \varepsilon_{\mathcal{T}}) - \mathcal{B}_{\mathcal{T}}(v_{\mathcal{T}}, \varepsilon_{\mathcal{T}})),$$

where we also used that $\mathcal{F}_{\mathcal{T}}(v) = (f, v)_{\Omega}$ since in this section we are working under a piecewise constant data assumption. We can split the above right hand side into two terms, obtaining

$$\|\varepsilon_{\mathcal{T}}\|^2 + |\varepsilon_{\mathcal{T}} - \mathcal{I}_{\mathcal{T}}\varepsilon_{\mathcal{T}}|_{1,\mathcal{T}}^2 \leq T_1 + T_2 \quad (8.6)$$

with

$$T_1 := \mathcal{B}(u - v_{\mathcal{T}}, \varepsilon_{\mathcal{T}}), \quad T_2 := \mathcal{B}(v_{\mathcal{T}}, \varepsilon_{\mathcal{T}}) - \mathcal{B}_{\mathcal{T}}(v_{\mathcal{T}}, \varepsilon_{\mathcal{T}}).$$

The bound for the first term is trivial

$$T_1 \leq \|u - v_{\mathcal{T}}\| \cdot \|\varepsilon_{\mathcal{T}}\|. \quad (8.7)$$

The second term is first written explicitly recalling the expression for $\mathcal{B}_{\mathcal{T}}(\cdot, \cdot)$, see (4.6), and using the orthogonality properties of the projectors

$$\begin{aligned} T_2 &= \sum_{E \in \mathcal{T}} \int_E (A_E \nabla(v_{\mathcal{T}} - \Pi_E^{\nabla} v_{\mathcal{T}})) \cdot \nabla \varepsilon_{\mathcal{T}} \\ &\quad + \sum_{E \in \mathcal{T}} \int_E (c_E(v_{\mathcal{T}} - \Pi_E^{\nabla} v_{\mathcal{T}})) \varepsilon_{\mathcal{T}} - s_E(v_{\mathcal{T}} - \mathcal{I}_E v_{\mathcal{T}}, \varepsilon_{\mathcal{T}} - \mathcal{I}_E \varepsilon_{\mathcal{T}}) \\ &\leq C(|v_{\mathcal{T}} - \Pi_E^{\nabla} v_{\mathcal{T}}|_{1,\mathcal{T}} + \|v_{\mathcal{T}} - \Pi_E^{\nabla} v_{\mathcal{T}}\|_{0,\Omega} + |v_{\mathcal{T}} - \mathcal{I}_{\mathcal{T}} v_{\mathcal{T}}|_{1,\mathcal{T}}) (\|\varepsilon_{\mathcal{T}}\| + |\varepsilon_{\mathcal{T}} - \mathcal{I}_{\mathcal{T}} \varepsilon_{\mathcal{T}}|_{1,\mathcal{T}}). \end{aligned}$$

Since the projector Π_E^{∇} minimizes the distance from (discontinuous) piecewise linear functions both in the broken H^1 semi-norm and in the L^2 norm, the above bound easily yields

$$T_2 \leq C|v_{\mathcal{T}} - \mathcal{I}_{\mathcal{T}} v_{\mathcal{T}}|_{1,\mathcal{T}} (\|\varepsilon_{\mathcal{T}}\| + |\varepsilon_{\mathcal{T}} - \mathcal{I}_{\mathcal{T}} \varepsilon_{\mathcal{T}}|_{1,\mathcal{T}}). \quad (8.8)$$

The result follows first combining bounds (8.6), (8.7), (8.8) and recalling (8.5). \square

Remark 8.2. Note that Lemma 8.1 would be false in the norm $\|\cdot\|$, that is without the second term in definition (8.2). Indeed this would imply that if $u \in \mathbb{V}_{\mathcal{T}}$ then $u_{\mathcal{T}} = u$, which is well known to be false in the VE method due to the approximation of the bilinear form.

We now introduce two different approximation classes, one based on the full Virtual Element space, and the other one based on the underlying piecewise linear conforming Finite Element space. Afterwards we will prove that, under the assumption of Λ -admissible partitions (cf. Definition 3.2), such classes are equivalent.

For any $N \in \mathbb{N}$, we define the following collection of partitions:

$$\mathbb{T}_N = \{\mathcal{T} : \mathcal{T} \text{ is } \Lambda\text{-admissible and satisfies } \#\mathcal{T} \leq N\}.$$

Definition 8.3 (approximation classes of v). *Given any $s \in \mathbb{R}$, $s > 0$, we define the following approximation classes*

$$\begin{aligned} \mathcal{A}_s &= \{v \in H_0^1(\Omega) : \exists C \in \mathbb{R} \text{ s.t. } \sigma_N(v) := \inf_{\mathcal{T} \in \mathbb{T}_N} \inf_{v_{\mathcal{T}} \in \mathbb{V}_{\mathcal{T}}} \mathbb{E}_{\mathcal{T}}(v, v_{\mathcal{T}}) \leq CN^{-s} \quad \forall N \geq \#\mathcal{T}_0\}, \\ \mathcal{A}_s^0 &= \{v \in H_0^1(\Omega) : \exists C \in \mathbb{R} \text{ s.t. } \sigma_N^0(v) := \inf_{\mathcal{T} \in \mathbb{T}_N} \inf_{v_{\mathcal{T}}^0 \in \mathbb{V}_{\mathcal{T}}^0} \mathbb{E}_{\mathcal{T}}(v, v_{\mathcal{T}}^0) \leq CN^{-s} \quad \forall N \geq \#\mathcal{T}_0\}. \end{aligned}$$

and denote

$$|v|_{\mathcal{A}_s} := \sup_{N \geq \#\mathcal{T}_0} N^s \sigma_N(v). \quad (8.9)$$

We now prove the following result on the equivalence of the approximation classes (see [12, Proposition 5.2]).

Proposition 8.4 (equivalence of classes). *The two classes in Definition 8.3 coincide, i.e.*

$$\mathcal{A}_s = \mathcal{A}_s^0 \quad \forall s \in \mathbb{R}, s > 0.$$

Proof. Let $s \in \mathbb{R}, s > 0$. The inclusion $\mathcal{A}_s^0 \subseteq \mathcal{A}_s$ is immediate since $\mathbb{V}_{\mathcal{T}}^0 \subseteq \mathbb{V}_{\mathcal{T}}$ and thus

$$\inf_{\mathcal{T} \in \mathbb{T}_N} \inf_{v_{\mathcal{T}} \in \mathbb{V}_{\mathcal{T}}} \mathbb{E}_{\mathcal{T}}^2(v, v_{\mathcal{T}}) \leq \inf_{\mathcal{T} \in \mathbb{T}_N} \inf_{v_{\mathcal{T}}^0 \in \mathbb{V}_{\mathcal{T}}^0} \mathbb{E}_{\mathcal{T}}^2(v, v_{\mathcal{T}}^0) \quad \forall v \in H_0^1(\Omega).$$

We now show the converse inclusion. We take a generic $v \in \mathcal{A}_s$. Let $N \geq \#\mathcal{T}_0$, then it exists $\mathcal{T} \in \mathbb{T}_N$ and $v_{\mathcal{T}} \in \mathbb{V}_{\mathcal{T}}$ such that

$$\mathbb{E}_{\mathcal{T}}^2(v, v_{\mathcal{T}}) = \|v - v_{\mathcal{T}}\|^2 + |v_{\mathcal{T}} - \mathcal{I}_{\mathcal{T}} v_{\mathcal{T}}|_{1,\mathcal{T}}^2 \leq CN^{-s},$$

with $C = C(v)$ but independent of N . We will exhibit an approximant in $\mathbb{V}_{\mathcal{T}}^0$ that satisfies the same bound, possibly with a different constant. We choose $\mathcal{I}_{\mathcal{T}}^0 v_{\mathcal{T}} \in \mathbb{V}_{\mathcal{T}}^0$, the Lagrange interpolant of $v_{\mathcal{T}}$ at the proper nodes of \mathcal{T} . Recalling observation (8.3) and by the triangle inequality

$$\begin{aligned} \mathbb{E}_{\mathcal{T}}^2(v, \mathcal{I}_{\mathcal{T}}^0 v_{\mathcal{T}}) &= \|v - \mathcal{I}_{\mathcal{T}}^0 v_{\mathcal{T}}\|^2 \leq 2(\|v - v_{\mathcal{T}}\|^2 + \|v_{\mathcal{T}} - \mathcal{I}_{\mathcal{T}}^0 v_{\mathcal{T}}\|^2) \\ &\leq C'(\|v - v_{\mathcal{T}}\|^2 + |v_{\mathcal{T}} - \mathcal{I}_{\mathcal{T}}^0 v_{\mathcal{T}}|_{1,\Omega}^2), \end{aligned}$$

where in the current proof C' denotes a generic positive constant that may change at each occurrence. Applying [8, Prop. 3.2] the above bound yields

$$\mathbb{E}_{\mathcal{T}}^2(v, \mathcal{I}_{\mathcal{T}}^0 v_{\mathcal{T}}) \leq C'(\|v - v_{\mathcal{T}}\|^2 + |v_{\mathcal{T}} - \mathcal{I}_{\mathcal{T}}^0 v_{\mathcal{T}}|_{1,\mathcal{T}}^2) \leq C' \mathbb{E}_{\mathcal{T}}^2(v, v_{\mathcal{T}}) \leq C' N^{-s}.$$

Since the constant C' does not depend on N , we have shown that $v \in \mathcal{A}_s^0$. Therefore $\mathcal{A}_s \subseteq \mathcal{A}_s^0$, and the proof is concluded. \square

In the rest of the paper, we make the following assumption.

Assumption 8.5 (approximability of u). *The solution u of Problem (2.1) belongs to \mathcal{A}_s for some $s = s_u \in (0, \frac{1}{2}]$.*

Remark 8.6 (equivalence with approximation classes on conforming partitions). It is easily seen that the class \mathcal{A}_s^0 , hence \mathcal{A}_s , coincides with the class \mathcal{A}_s^c defined by replacing \mathbb{T}_N by $\mathbb{T}_N^c = \{\mathcal{T} : \mathcal{T} \text{ is conforming and satisfies } \#\mathcal{T} \leq N\}$. Indeed, any $\mathcal{T} \in \mathbb{T}_N$ can be refined to produce a conforming partition \mathcal{T}^c , such that $\#\mathcal{T}^c \leq K \#\mathcal{T}$ for a positive constant $K = K_{\Lambda}$ solely depending on Λ . As a consequence, one can apply e.g. [10, Theorem 9.1] and deduce that $u \in \mathcal{A}_{\frac{1}{2}}^c$ provided $u \in W_p^2(\Omega)$ for some $p > 1$.

It must be finally observed that the important result above does not exclude that AVEM, which contains AFEM and allows more flexibility in terms of hanging nodes, could obtain a better efficiency in terms of the involved constants (in this respect, see also Section 10).

8.1.2 Approximation classes for data

Given a partition \mathcal{T} and piecewise constant data $\widehat{\mathcal{D}} = (\widehat{A}_{\mathcal{T}}, \widehat{c}_{\mathcal{T}}, \widehat{f}_{\mathcal{T}})$ defined as in (6.1), let us set (cf. (6.11))

$$\zeta_{\mathcal{T}}(A) = \|A - \widehat{A}_{\mathcal{T}}\|_{L^\infty(\Omega)}, \quad \zeta_{\mathcal{T}}(c) = \|h(c - \widehat{c}_{\mathcal{T}})\|_{L^\infty(\Omega)}, \quad \zeta_{\mathcal{T}}(f) = \|h(f - \widehat{f}_{\mathcal{T}})\|_{L^2(\Omega)}. \quad (8.10)$$

Definition 8.7 (approximation classes of A). *Let*

$$\mathbb{A}_s = \{A \in (L^\infty(\Omega))^{2 \times 2} : \exists C \in \mathbb{R} \text{ s.t. } \inf_{\mathcal{T} \in \mathbb{T}_N} \zeta_{\mathcal{T}}(A) \leq CN^{-s} \ \forall N \geq \#\mathcal{T}_0\} \quad (8.11)$$

and denote

$$|A|_{\mathbb{A}_s} := \sup_{N \geq \#\mathcal{T}_0} \left(N^s \inf_{\mathcal{T} \in \mathbb{T}_N} \zeta_{\mathcal{T}}(A) \right). \quad (8.12)$$

Definition 8.8 (approximation classes of c). *Let*

$$\mathbb{C}_s = \{c \in L^\infty(\Omega) : \exists C \in \mathbb{R} \text{ s.t. } \inf_{\mathcal{T} \in \mathbb{T}_N} \zeta_{\mathcal{T}}(c) \leq CN^{-s} \ \forall N \geq \#\mathcal{T}_0\} \quad (8.13)$$

and denote

$$|c|_{\mathbb{C}_s} := \sup_{N \geq \#\mathcal{T}_0} \left(N^s \inf_{\mathcal{T} \in \mathbb{T}_N} \zeta_{\mathcal{T}}(c) \right). \quad (8.14)$$

Definition 8.9 (approximation classes of f). *Let*

$$\mathbb{F}_s = \{f \in L^2(\Omega) : \exists C \in \mathbb{R} \text{ s.t. } \inf_{\mathcal{T} \in \mathbb{T}_N} \zeta_{\mathcal{T}}(f) \leq CN^{-s} \ N \geq \#\mathcal{T}_0\} \quad (8.15)$$

and denote

$$|f|_{\mathbb{F}_s} := \sup_{N \geq \#\mathcal{T}_0} \left(N^s \inf_{\mathcal{T} \in \mathbb{T}_N} \zeta_{\mathcal{T}}(f) \right). \quad (8.16)$$

In the rest of the paper, we make the following assumptions concerning the data of our problem and their piecewise-linear approximation.

Assumption 8.10 (approximability of data). *There exist $s_A, s_c, s_f \in (0, \frac{1}{2}]$ such that the data of Problem (2.1) satisfy $A \in \mathbb{A}_{s_A}$, $c \in \mathbb{C}_{s_c}$, $f \in \mathbb{F}_{s_f}$.*

Assumption 8.11 (quasi-optimality of the module DATA). *The procedure MARK_DATA introduced in Sect. 6.2 is quasi-optimal, namely the cardinalities of the marked sets $\mathcal{M}_A, \mathcal{M}_c, \mathcal{M}_f$ for A, c, f resp., satisfy*

$$\#\mathcal{M}_A \lesssim |A|_{\mathbb{A}_{s_A}}^{\frac{1}{s_A}} \varepsilon^{-\frac{1}{s_A}}, \quad \#\mathcal{M}_c \lesssim |c|_{\mathbb{C}_{s_c}}^{\frac{1}{s_c}} \varepsilon^{-\frac{1}{s_c}}, \quad \#\mathcal{M}_f \lesssim |f|_{\mathbb{F}_{s_f}}^{\frac{1}{s_f}} \varepsilon^{-\frac{1}{s_f}}. \quad (8.17)$$

Under this assumption, setting $s_{\mathcal{D}} = \min(s_A, s_c, s_f)$, the cardinality of the marked set $\mathcal{M}_{\mathcal{D}} = \mathcal{M}_A \cup \mathcal{M}_c \cup \mathcal{M}_f$ satisfies

$$\#\mathcal{M}_{\mathcal{D}} \lesssim (|A|_{\mathbb{A}_{s_A}}^{\frac{1}{s_A}} + |c|_{\mathbb{C}_{s_c}}^{\frac{1}{s_c}} + |f|_{\mathbb{F}_{s_f}}^{\frac{1}{s_f}}) \varepsilon^{-\frac{1}{s_{\mathcal{D}}}} =: |\mathcal{D}|_{\mathbb{A}_{\mathcal{D}}}^{\frac{1}{s_{\mathcal{D}}}} \varepsilon^{-\frac{1}{s_{\mathcal{D}}}}. \quad (8.18)$$

Remark 8.12. In Sect. 9 we will give regularity conditions on the data such that Assumption 8.10 is satisfied. In particular, we will prove that $s_A = s_c = s_f = \frac{1}{2}$ if $A \in (W_p^1(\Omega))^{2 \times 2}$ with $p > 1$, $c \in L^\infty(\Omega)$ and $f \in L^2(\Omega)$. Furthermore, we will show that the implementation of MARK_DATA described in Sect. 6.2 guarantees the validity of Assumption 8.11.

8.2 ε -approximation of order s

Since the data $\widehat{\mathcal{D}}_k$ is fixed inside GALERKIN, the performance of this module is dictated by the regularity of $\widehat{u}_k = \text{exact.sol}(\widehat{\mathcal{D}}_k)$, which is the exact solution with data $\widehat{\mathcal{D}}_k$, rather than u . We know that $u \in \mathcal{A}_s$ and wonder what regularity is inherited by \widehat{u}_k . This leads to the following concept introduced in [11, Def. 3.1 and Lemma 3.2].

Definition 8.13 (ε -approximation of order s). *Given $u \in \mathcal{A}_s$ and $\varepsilon > 0$, a function $v \in H_0^1(\Omega)$ is said to be an ε -approximation of order s to u if $\|u - v\| \leq \varepsilon$ and there exists a constant $C > 0$ independent of ε , u and v such that for all $\delta \geq \varepsilon$ there exists $N \geq \#\mathcal{T}_0$ satisfying*

$$\sigma_N(v) \leq \delta \quad N \leq C|u|_{\mathcal{A}_s}^{\frac{1}{s}} \delta^{-\frac{1}{s}} + 1.$$

Remark 8.14. In view of the definition of $\sigma_N(v)$, there exists $\mathcal{T} \in \mathbb{T}_N$ and $v_{\mathcal{T}} \in \mathbb{V}_{\mathcal{T}}$ such that

$$\sigma_N(v) = \mathbb{E}_{\mathcal{T}}(v, v_{\mathcal{T}}) \leq \delta.$$

Lemma 8.15 (ε -approximation of u of order s). *Let $u \in \mathcal{A}_s$ and $v \in H_0^1(\Omega)$ satisfying $\|u - v\| \leq \varepsilon$ for some $\varepsilon > 0$. Then v is a 2ε -approximation of order s to u .*

Proof. Let $\delta \geq 2\varepsilon$. By definition of $\sigma_N(u)$, there exists $N \geq \#\mathcal{T}_0$, $\mathcal{T} \in \mathbb{T}_N$ and $w_{\mathcal{T}} \in \mathbb{V}_{\mathcal{T}}$ such that

$$\sigma_N(u) = \mathbb{E}_{\mathcal{T}}(u, w_{\mathcal{T}}) \leq \frac{\delta}{4} \quad N \leq |u|_{\mathcal{A}_s}^{\frac{1}{s}} \left(\frac{\delta}{4}\right)^{-\frac{1}{s}} + 1.$$

The triangle and Young inequalities yield

$$\begin{aligned} \sigma_N(v) &\leq \mathbb{E}_{\mathcal{T}}(v, w_{\mathcal{T}}) \leq \|v - w_{\mathcal{T}}\| + |w_{\mathcal{T}} - \mathcal{I}_{\mathcal{T}} w_{\mathcal{T}}|_{1,\Omega} \\ &\leq \|v - u\| + \|u - w_{\mathcal{T}}\| + |w_{\mathcal{T}} - \mathcal{I}_{\mathcal{T}} w_{\mathcal{T}}|_{1,\Omega} \leq \|v - u\| + \sqrt{2} \mathbb{E}_{\mathcal{T}}(u, w_{\mathcal{T}}) \\ &\leq \varepsilon + \sqrt{2} \frac{\delta}{4} \leq \left(\frac{1}{2} + \frac{\sqrt{2}}{2}\right) \delta < \delta. \end{aligned}$$

Moreover, there holds

$$N \leq 4^{\frac{1}{s}} |u|_{\mathcal{A}_s}^{\frac{1}{s}} \delta^{-\frac{1}{s}} + 1.$$

This concludes the proof with constant $C = 4^{\frac{1}{s}}$. \square

8.3 Optimality of mesh refinement

Hereafter, we consider two Λ -admissible partitions \mathcal{T} and \mathcal{T}_* , the latter being a refinement of the former obtained by applying a newest-vertex bisection to some of the elements of \mathcal{T} . Considering the corresponding Galerkin solutions $u_{\mathcal{T}}$ and $u_{\mathcal{T}_*}$ of problem (4.10) with piecewise constant data, we first prove that the difference in energy norm between $u_{\mathcal{T}}$ and the orthogonal projection of $u_{\mathcal{T}_*}$ upon $\mathbb{V}_{\mathcal{T}_*}^0$ can be essentially bounded by the contribution to the error estimator coming from a neighborhood of the refined elements. Next, we give conditions under which this portion of the error estimator satisfies a Dörfler property with respect to the full estimator.

8.3.1 Localized upper bound of the difference between Galerkin solutions

Consider an element $E \in \mathcal{T}$ which has been split into two elements $E_1, E_2 \in \mathcal{T}_*$. If $v \in \mathbb{V}_{\mathcal{T}}$, then v is known on ∂E , hence in particular at the new vertex of E_1, E_2 created by bisection. Thus, v is known at all nodes (vertices and possibly hanging nodes) sitting on ∂E_1 and ∂E_2 , since the new edge $e = E_1 \cap E_2$ does not contain internal nodes. This uniquely identifies a function in \mathbb{V}_{E_1} and a function in \mathbb{V}_{E_2} , which are continuous across e . In this manner, we associate to any $v \in \mathbb{V}_{\mathcal{T}}$ a unique function $v_* \in \mathbb{V}_{\mathcal{T}_*}$, that coincides with v on the skeleton \mathcal{E} . We will actually write v for v_* whenever no confusion is possible.

We introduce the following orthogonal decomposition of $\mathbb{V}_{\mathcal{T}}$

$$\mathbb{V}_{\mathcal{T}} = \mathbb{V}_{\mathcal{T}}^0 \oplus \mathbb{V}_{\mathcal{T}}^{\perp}, \quad (8.19)$$

where $\mathbb{V}_{\mathcal{T}}^{\perp}$ is the orthogonal complement of $\mathbb{V}_{\mathcal{T}}^0$ in $\mathbb{V}_{\mathcal{T}}$ with respect to the (discrete) scalar product $\mathcal{B}_{\mathcal{T}}(\cdot, \cdot)$, and we prove a localized estimate (cf. [12, Lemma 3.5]) that is crucial in the discussion of the quasi-optimal cardinality of our adaptive algorithm. To this end, we denote by $\mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*}$ the set of refined elements of \mathcal{T} to obtain \mathcal{T}_* and let $\omega(\mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*})$ be any subset of \mathcal{T} containing $\mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*}$. We observe that as \mathcal{T}_* is a refinement of \mathcal{T} , Assumption 4.1 implies that for every $E_* \in \mathcal{T}_*$ with $E_* \subseteq E$, $E \in \mathcal{T}$ we have $A_{E_*} = A_E$, $c_{E_*} = c_E$ and $f_{E_*} = f_E$. The following lemma bounds the difference between a discrete solution and (the $\mathbb{V}_{\mathcal{T}_*}^0$ part of) another discrete solution on a refined mesh. Such difference is bounded by the error estimator evaluated on a suitable neighbourhood of the refined elements, plus an additional term which nevertheless becomes “negligible” for γ sufficiently large.

Lemma 8.16 (localized upper bound). *Let \mathcal{T}_* be a refinement of \mathcal{T} and let $u_{\mathcal{T}} \in \mathbb{V}_{\mathcal{T}}$ and $u_{\mathcal{T}_*} \in \mathbb{V}_{\mathcal{T}_*}$ be the corresponding discrete solutions of (4.10) with piecewise constant data. Let $u_{\mathcal{T}_*} = u_{\mathcal{T}_*}^0 + u_{\mathcal{T}_*}^{\perp} \in \mathbb{V}_{\mathcal{T}_*}^0 \oplus \mathbb{V}_{\mathcal{T}_*}^{\perp}$ be the orthogonal decomposition of $u_{\mathcal{T}_*}$ according to (8.19). Then, there exists a constant C_{LU} only depending on the shape regularity of \mathcal{T} so that*

$$\|u_{\mathcal{T}_*}^0 - u_{\mathcal{T}}\| \leq C_{LU} (\eta_{\mathcal{T}}(\omega(\mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*}); u_{\mathcal{T}}, \mathcal{D}) + \gamma^{-1} \eta_{\mathcal{T}}(u_{\mathcal{T}}, \mathcal{D})) . \quad (8.20)$$

Proof. Let us preliminarily proceed by steps and collect some instrumental results that will be employed in the sequel.

Step 1. First, we observe that as $\mathbb{V}_{\mathcal{T}}^0 \subset \mathbb{V}_{\mathcal{T}_*}^0$ is made of continuous piecewise linear functions on \mathcal{T} we have

$$v_{\mathcal{T}}^0 = \Pi_{\mathcal{T}}^{\nabla} v_{\mathcal{T}}^0 = \Pi_{\mathcal{T}_*}^{\nabla} v_{\mathcal{T}}^0. \quad (8.21)$$

Step 2. There holds

$$\mathcal{B}_{\mathcal{T}_*}(u_{\mathcal{T}_*}^0, v_{\mathcal{T}_*}^0) = \mathcal{F}_{\mathcal{T}_*}(v_{\mathcal{T}_*}^0) \quad \forall v_{\mathcal{T}_*}^0 \in \mathbb{V}_{\mathcal{T}_*}^0. \quad (8.22)$$

Indeed, for any $v_{\mathcal{T}_*}^0 \in \mathbb{V}_{\mathcal{T}_*}^0$ we have

$$\mathcal{F}_{\mathcal{T}_*}(v_{\mathcal{T}_*}^0) = \mathcal{B}_{\mathcal{T}_*}(u_{\mathcal{T}_*}, v_{\mathcal{T}_*}^0) = \mathcal{B}_{\mathcal{T}_*}(u_{\mathcal{T}_*}^0 + u_{\mathcal{T}_*}^{\perp}, v_{\mathcal{T}_*}^0) = \mathcal{B}_{\mathcal{T}_*}(u_{\mathcal{T}_*}^0, v_{\mathcal{T}_*}^0) \quad (8.23)$$

where in the last step we employed that $\mathbb{V}_{\mathcal{T}_*}^{\perp}$ is the orthogonal complement of $\mathbb{V}_{\mathcal{T}_*}^0$ in $\mathbb{V}_{\mathcal{T}_*}$ with respect to $\mathcal{B}_{\mathcal{T}_*}(\cdot, \cdot)$.

Step 3. There holds

$$\mathcal{B}(u_{\mathcal{T}_*}^0 - u_{\mathcal{T}}, v_{\mathcal{T}}^0) = 0 \quad \forall v_{\mathcal{T}}^0 \in \mathbb{V}_{\mathcal{T}}^0. \quad (8.24)$$

Using (4.9) and (8.21) we have

$$\mathcal{B}(u_{\mathcal{T}_*}^0 - u_{\mathcal{T}}, v_{\mathcal{T}}^0) = \mathcal{B}_{\mathcal{T}_*}(u_{\mathcal{T}_*}^0, v_{\mathcal{T}}^0) - \mathcal{B}_{\mathcal{T}}(u_{\mathcal{T}}, v_{\mathcal{T}}^0) = \mathcal{F}_{\mathcal{T}_*}(v_{\mathcal{T}}^0) - \mathcal{F}_{\mathcal{T}}(v_{\mathcal{T}}^0) \quad (8.25)$$

where in the last step we employed (8.22). From Assumption 4.1, (4.11) and (8.21) we get (8.24).

Step 4. Let $e_*^0 = u_{\mathcal{T}_*}^0 - u_{\mathcal{T}} - v_{\mathcal{T}}^0$ with $v_{\mathcal{T}}^0 \in V_{\mathcal{T}}^0$, where $u_{\mathcal{T}}^0 = u_{\mathcal{T}} - u_{\mathcal{T}}^\perp \in \mathbb{V}_{\mathcal{T}}^0$. There holds

$$\mathcal{B}_{\mathcal{T}_*}(u_{\mathcal{T}_*}^0 - u_{\mathcal{T}}, e_*^0) \lesssim |u_{\mathcal{T}_*}^0 - u_{\mathcal{T}}^0|_{1,\Omega} (\eta_{\mathcal{T}}(\omega(\mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*}); u_{\mathcal{T}}, \mathcal{D}) + \gamma^{-1} \eta_{\mathcal{T}}(u_{\mathcal{T}}, \mathcal{D})). \quad (8.26)$$

Indeed, we have

$$\begin{aligned} \mathcal{B}_{\mathcal{T}_*}(u_{\mathcal{T}_*}^0 - u_{\mathcal{T}}, e_*^0) &= \mathcal{F}_{\mathcal{T}_*}(e_*^0) - \mathcal{B}_{\mathcal{T}_*}(u_{\mathcal{T}}, e_*^0) \\ &= (\mathcal{F}_{\mathcal{T}_*}(e_*^0) - \mathcal{B}_{\mathcal{T}_*}(\Pi_{\mathcal{T}}^\nabla u_{\mathcal{T}}, e_*^0)) + \mathcal{B}_{\mathcal{T}_*}(\Pi_{\mathcal{T}}^\nabla u_{\mathcal{T}} - u_{\mathcal{T}}, e_*^0) =: I + II, \end{aligned}$$

where, with a slight abuse of notation, we extend the definition (4.6) of $\mathcal{B}_{\mathcal{T}_*}$ to $\mathbb{P}_1(\mathcal{T}_*)$.

In the sequel we choose $v_{\mathcal{T}}^0 = \tilde{\mathcal{I}}_{\mathcal{T}}^0(u_{\mathcal{T}_*}^0 - u_{\mathcal{T}}^0)$ in the definition of e_*^0 , where $\tilde{\mathcal{I}}_{\mathcal{T}}^0 : C^0(\bar{\Omega}) \rightarrow \mathbb{V}_{\mathcal{T}}^0$ is the Clément quasi-interpolation operator on \mathcal{T}^0 . We also notice that e_*^0 vanishes outside the set $\omega(\mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*})$. As $e_*^0 \in V_{\mathcal{T}_*}^0$ and $\Pi_{\mathcal{T}}^\nabla e_*^0 = e_*^0$ we have

$$\begin{aligned} I &= \mathcal{F}_{\mathcal{T}_*}(e_*^0) - \sum_{E_* \in \mathcal{T}_*} \int_{E_*} (A_{E_*} \nabla \Pi_{\mathcal{T}_*}^\nabla (\Pi_{\mathcal{T}}^\nabla u_{\mathcal{T}}) \cdot \nabla \Pi_{\mathcal{T}_*}^\nabla e_*^0 + c_{E_*} \Pi_{\mathcal{T}_*}^\nabla (\Pi_{\mathcal{T}}^\nabla u_{\mathcal{T}}) \Pi_{\mathcal{T}_*}^\nabla e_*^0) \\ &= \mathcal{F}_{\mathcal{T}}(e_*^0) - \sum_{E \in \mathcal{T}} \sum_{E_* \in \mathcal{T}_*, E_* \subseteq E} \int_{E_*} (A_{E_*} \nabla \Pi_{\mathcal{T}}^\nabla u_{\mathcal{T}} \cdot \nabla e_*^0 + c_{E_*} \Pi_{\mathcal{T}}^\nabla u_{\mathcal{T}} e_*^0) \\ &= \sum_{E \in \omega(\mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*})} \int_E (f_E e_*^0 - A_E \nabla \Pi_{\mathcal{T}}^\nabla u_{\mathcal{T}} \cdot \nabla e_*^0 - c_E \Pi_{\mathcal{T}}^\nabla u_{\mathcal{T}} e_*^0) \end{aligned}$$

where we employed the properties of the enhanced space (3.2). Integrating by parts, employing the Cauchy-Schwarz inequality together with the vanishing property of e_*^0 and the interpolation error estimate for $\tilde{\mathcal{I}}_{\mathcal{T}}^0$, we obtain

$$I \lesssim |u_{\mathcal{T}_*}^0 - u_{\mathcal{T}}^0|_{1,\Omega} \eta_{\mathcal{T}}(\omega(\mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*}); u_{\mathcal{T}}, \mathcal{D}). \quad (8.27)$$

On the other hand, again as $e_*^0 \in V_{\mathcal{T}_*}^0$ and $\Pi_{\mathcal{T}}^\nabla e_*^0 = e_*^0$, we have

$$\begin{aligned} II &= \sum_{E_* \in \mathcal{T}_*} \int_{E_*} (A_{E_*} \nabla \Pi_{\mathcal{T}_*}^\nabla (\Pi_{\mathcal{T}}^\nabla - I) u_{\mathcal{T}} \cdot \nabla \Pi_{\mathcal{T}_*}^\nabla e_*^0 + c_{E_*} \Pi_{\mathcal{T}_*}^\nabla (\Pi_{\mathcal{T}}^\nabla - I) u_{\mathcal{T}} \Pi_{\mathcal{T}_*}^\nabla e_*^0) \\ &= \sum_{E_* \in \mathcal{T}_*} \int_{E_*} (A_{E_*} \nabla (\Pi_{\mathcal{T}}^\nabla - I) u_{\mathcal{T}} \cdot \nabla e_*^0 + c_{E_*} (\Pi_{\mathcal{T}}^\nabla - I) u_{\mathcal{T}} e_*^0) \\ &= \sum_{E \in \mathcal{T}} \sum_{E_* \in \mathcal{T}_*, E_* \subseteq E} \int_{E_*} (A_{E_*} \nabla (\Pi_{\mathcal{T}}^\nabla - I) u_{\mathcal{T}} \cdot \nabla e_*^0 + c_{E_*} (\Pi_{\mathcal{T}}^\nabla - I) u_{\mathcal{T}} e_*^0) \\ &= \sum_{E \in \mathcal{T}} \int_E (A_E \nabla (\Pi_{\mathcal{T}}^\nabla - I) u_{\mathcal{T}} \cdot \nabla e_*^0 + c_E (\Pi_{\mathcal{T}}^\nabla - I) u_{\mathcal{T}} e_*^0) \\ &\lesssim S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}})^{1/2} |u_{\mathcal{T}_*}^0 - u_{\mathcal{T}}^0|_{1,\Omega} \lesssim \gamma^{-1} \eta_{\mathcal{T}}(u_{\mathcal{T}}, \mathcal{D}) |u_{\mathcal{T}_*}^0 - u_{\mathcal{T}}^0|_{1,\Omega}. \end{aligned} \quad (8.28)$$

where in the last step we used (4.19). The thesis follows combining (8.27)-(8.28).

Step 5. Let $u_{\mathcal{T}} = u_{\mathcal{T}}^0 + u_{\mathcal{T}}^{\perp}$ be the orthogonal decomposition (8.19). There holds

$$\|u_{\mathcal{T}}^{\perp}\| \lesssim S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}})^{1/2} \lesssim \gamma^{-1} \eta_{\mathcal{T}}(u_{\mathcal{T}}, \mathcal{D}). \quad (8.29)$$

Indeed, we have

$$\mathcal{B}_{\mathcal{T}}(u_{\mathcal{T}}^{\perp}, u_{\mathcal{T}}^{\perp}) = \inf_{w_{\mathcal{T}}^0 \in V_{\mathcal{T}}^0} \mathcal{B}_{\mathcal{T}}(u_{\mathcal{T}} - w_{\mathcal{T}}^0, u_{\mathcal{T}} - w_{\mathcal{T}}^0) \leq \mathcal{B}_{\mathcal{T}}(u_{\mathcal{T}} - I_{\mathcal{T}}^0 u_{\mathcal{T}}, u_{\mathcal{T}} - I_{\mathcal{T}}^0 u_{\mathcal{T}}) \lesssim S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}})$$

where in the last inequality we employed the continuity of $\mathcal{B}_{\mathcal{T}}(\cdot, \cdot)$ in combination with [8, Prop. 3.2]. The coercivity of $\mathcal{B}_{\mathcal{T}}(\cdot, \cdot)$ together with (4.5) and (2.3), and the bound (4.19) yield the result.

At this point, we have collected all ingredients to prove (8.20). From the coercivity of $\mathcal{B}(\cdot, \cdot)$ and employing (8.24) we get

$$\begin{aligned} \|u_{\mathcal{T}_*}^0 - u_{\mathcal{T}}\|^2 &\lesssim \mathcal{B}(u_{\mathcal{T}_*}^0 - u_{\mathcal{T}}, u_{\mathcal{T}_*}^0 - u_{\mathcal{T}}) = \mathcal{B}(u_{\mathcal{T}_*}^0 - u_{\mathcal{T}}, e_*) + \mathcal{B}(u_{\mathcal{T}_*}^0 - u_{\mathcal{T}}, v_{\mathcal{T}}^0 - u_{\mathcal{T}}^{\perp}) \\ &= \mathcal{B}(u_{\mathcal{T}_*}^0 - u_{\mathcal{T}}, e_*) - \mathcal{B}(u_{\mathcal{T}_*}^0 - u_{\mathcal{T}}, u_{\mathcal{T}}^{\perp}) =: III + IV. \end{aligned}$$

Employing the consistency of $\mathcal{B}_{\mathcal{T}_*}(\cdot, \cdot)$ (cf. (4.9)) together with (8.26) we get

$$III = \mathcal{B}_{\mathcal{T}_*}(u_{\mathcal{T}_*}^0 - u_{\mathcal{T}}, e_*) \lesssim |u_{\mathcal{T}_*}^0 - u_{\mathcal{T}}|_{1,\Omega} (\eta_{\mathcal{T}}(\omega(\mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*}); u_{\mathcal{T}}, \mathcal{D}) + \gamma^{-1} \eta_{\mathcal{T}}(u_{\mathcal{T}}, \mathcal{D})). \quad (8.30)$$

On the other hand, employing the continuity of $\mathcal{B}(\cdot, \cdot)$ in combination with (8.29) we obtain

$$IV \lesssim \|u_{\mathcal{T}_*}^0 - u_{\mathcal{T}}\| \|u_{\mathcal{T}}^{\perp}\| \lesssim \|u_{\mathcal{T}_*}^0 - u_{\mathcal{T}}\| \gamma^{-1} \eta_{\mathcal{T}}(u_{\mathcal{T}}, \mathcal{D}). \quad (8.31)$$

We now observe that $|u_{\mathcal{T}_*}^0 - u_{\mathcal{T}}|_{1,\Omega} = |u_{\mathcal{T}_*}^0 - u_{\mathcal{T}} + u_{\mathcal{T}}^{\perp}|_{1,\Omega} \lesssim \|u_{\mathcal{T}_*}^0 - u_{\mathcal{T}}\| + \gamma^{-1} \eta_{\mathcal{T}}(u_{\mathcal{T}}, \mathcal{D})$. Concatenating (8.30)-(8.31), we easily conclude the proof of Lemma 8.16. \square

8.3.2 Optimal marking

We first recall two instrumental results that will be useful in the sequel. From [8, Corollary 4.3] we have the global error bound

$$C_{GL} \eta_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}) \leq \|u - u_{\mathcal{T}}\|^2 + |u_{\mathcal{T}} - I_{\mathcal{T}} u_{\mathcal{T}}|_{1,\mathcal{T}}^2 = \mathbb{E}_{\mathcal{T}}^2(u, u_{\mathcal{T}}). \quad (8.32)$$

Moreover, we observe that (4.5) and (4.19) yield

$$|u_{\mathcal{T}} - I_{\mathcal{T}} u_{\mathcal{T}}|_{1,\mathcal{T}}^2 \leq \tilde{C}_B \gamma^{-2} \eta_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}). \quad (8.33)$$

In order to derive a quasi-optimal decay of the total error, we define

$$\gamma_*^2 := \frac{2C_{LU} + \tilde{C}_B}{C_{GL}} \quad \theta_*(\gamma) := \frac{C_{GL} - \gamma^{-2}(2C_{LU} + \tilde{C}_B)}{2C_{LU}}$$

for $\gamma > \gamma_*$, where C_{LU} is given by Lemma 8.16. Notice that $\gamma > \gamma_*$ yields $\theta_* > 0$ and if $C_{GL} < 2C_{LU}$ then $\theta_* < 1$. Moreover we make the following assumption.

Assumption 8.17 (module MARK). *The set of marked elements produced by the module MARK has minimal cardinality and the marking parameter satisfies $\theta \in (0, \theta_*)$.*

In order to simplify the notation, we let $0 < \mu < 1/2$ be defined by

$$\mu(\gamma, \theta) := \frac{C_{GL} - \gamma^{-2}(2C_{LU} + \tilde{C}_B)}{2C_{GL}} \left(1 - \frac{\theta}{\theta_*}\right) \quad \forall \gamma > \gamma_*, \quad 0 < \theta < \theta_*.$$

We now prove the analogous of [12, Lemma 5.4].

Lemma 8.18 (optimal marking). *Let \mathcal{T}_* be a refinement of \mathcal{T} and let $u_{\mathcal{T}} \in V_{\mathcal{T}}$ and $u_{\mathcal{T}_*} \in V_{\mathcal{T}_*}$ the corresponding discrete solutions of (4.10). In addition, assume*

$$\mathbb{E}_{\mathcal{T}}^2(u, u_{\mathcal{T}_*}^0) \leq \mu \mathbb{E}_{\mathcal{T}}^2(u, u_{\mathcal{T}}) \quad (8.34)$$

where $u_{\mathcal{T}_*} = u_{\mathcal{T}_*}^0 + u_{\mathcal{T}_*}^\perp$ is the orthogonal decomposition (8.19). Then, for $\gamma > \gamma_*$ and $\theta \in (0, \theta_*(\gamma))$, the set $\omega(\mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*})$ satisfies a Dörfler marking property

$$\eta_{\mathcal{T}}^2(\omega(\mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*}); u_{\mathcal{T}}, \mathcal{D}) \geq \theta \eta_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}).$$

Proof. Since $0 < \mu < 1/2$, employing (8.32) and (8.34) we get

$$(1 - 2\mu)C_{GL}\eta_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}) \leq \|u - u_{\mathcal{T}}\|^2 - 2\|u - u_{\mathcal{T}_*}^0\|^2 + |u_{\mathcal{T}} - I_{\mathcal{T}}u_{\mathcal{T}}|_{1,\mathcal{T}}^2 \quad (8.35)$$

where we used $I_{\mathcal{T}_*}u_{\mathcal{T}_*}^0 = u_{\mathcal{T}_*}^0$. From the triangle inequality and (8.20) we obtain

$$\|u - u_{\mathcal{T}}\|^2 - 2\|u - u_{\mathcal{T}_*}^0\|^2 \leq 2\|u_{\mathcal{T}} - u_{\mathcal{T}_*}^0\|^2 \leq 2C_{LU} \left(\eta_{\mathcal{T}}^2(\omega(\mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*}); u_{\mathcal{T}}, \mathcal{D}) + \gamma^{-2} \eta_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}) \right). \quad (8.36)$$

Combining (8.35)-(8.33) we get

$$(1 - 2\mu)C_{GL}\eta_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}) \leq 2C_{LU}\eta_{\mathcal{T}}^2(\omega(\mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*}); u_{\mathcal{T}}, \mathcal{D}) + \gamma^{-2}(2C_{LU} + \tilde{C}_B)\eta_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D})$$

which implies, employing the definition of μ and θ_* , the desired estimate

$$\eta_{\mathcal{T}}^2(\omega(\mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*}); u_{\mathcal{T}}, \mathcal{D}) \geq \frac{1}{2C_{LU}} \left((1 - 2\mu)C_{GL} - \gamma^{-2}(2C_{LU} + \tilde{C}_B) \right) \eta_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}) \geq \theta \eta_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}).$$

The proof is concluded. \square

8.4 Complexity of GALERKIN

In this section, we rely on the notation introduced in Sect. 7, in particular those in (7.9). We assume that the pair $(\hat{\mathcal{T}}_k, \hat{\mathcal{D}}_k)$ transferred by DATA to GALERKIN at iteration k satisfies

$$\eta_{\hat{\mathcal{T}}_k}(u_{k,0}, \hat{\mathcal{D}}_k) =: \hat{\varepsilon}_k > \varepsilon_k,$$

for otherwise GALERKIN is skipped. On the other hand, combining (2.3) with the stabilization free a posteriori error estimates (4.20), we can write

$$\hat{C}_L^2 \eta_{\hat{\mathcal{T}}_k}^2(u_{k,0}, \hat{\mathcal{D}}_k) \leq \|\hat{u}_k - u_{k,0}\|^2 \leq \hat{C}_U^2 \eta_{\hat{\mathcal{T}}_k}^2(u_{k,0}, \hat{\mathcal{D}}_k) \quad (8.37)$$

with $\hat{C}_L^2 := c_{\mathcal{B}}C_L$ and $\hat{C}_U^2 := c^{\mathcal{B}}C_U$. Therefore, we get the lower bound

$$\|\hat{u}_k - u_{k,0}\| \geq \hat{C}_L \hat{\varepsilon}_k > \hat{C}_L \varepsilon_k.$$

On the other hand, from (6.7) and (6.9) it follows that DATA provides a perturbed exact solution $\hat{u}_k \in H_0^1(\Omega)$ satisfying

$$\|u - \hat{u}_k\| \leq D\omega\varepsilon_k = \frac{D\omega}{\hat{C}_L} \hat{C}_L \varepsilon_k$$

for a suitable constant $D > 0$. Let

$$\omega := \frac{\sqrt{\mu}\hat{C}_L}{2D},$$

which implies

$$\|u - \hat{u}_k\| \leq \frac{\sqrt{\mu}}{2} \hat{C}_L \varepsilon_k.$$

In view of Proposition 7.3 (computational cost of GALERKIN) the module GALERKIN performs a number of iterations J_k bounded uniformly in k by J . For each such iteration j we have a mesh $\mathcal{T}_{k,j}$ and a Galerkin solution $u_{k,j} \in \mathbb{V}_{\mathcal{T}_{k,j}}$ so that for $0 \leq j < J_k$

$$\begin{aligned} \mathcal{T}_{k,0} &= \hat{\mathcal{T}}_k, \\ \eta_{\mathcal{T}_{k,j}}(u_{k,j}, \hat{\mathcal{D}}_k) &> \varepsilon_k, \\ \mathbb{E}_{\mathcal{T}_{k,j}}(\hat{u}_k, u_{k,j}) &\geq \|\hat{u}_k - u_{k,j}\| \geq \hat{C}_L \eta_{\mathcal{T}_{k,j}}(u_{k,j}, \hat{\mathcal{D}}_k) > \hat{C}_L \varepsilon_k. \end{aligned}$$

Let $\mathcal{M}_{k,j}$ be the marked set within $\mathcal{T}_{k,j}$ using the Dörfler strategy.

Lemma 8.19 (cardinality of marked sets). *If $u \in \mathcal{A}_s$ and $\omega = \frac{\sqrt{\mu}\hat{C}_L}{2D}$, then there exists a constant $C_0 > 0$ such that*

$$\#\mathcal{M}_{k,j} \leq C_0 |u|_{\mathcal{A}_s}^{\frac{1}{s}} \varepsilon^{-\frac{1}{s}}, \quad 0 \leq j < J_k.$$

Proof. Fix $0 \leq j < J_k$ and set

$$\delta := \sqrt{\mu} \mathbb{E}_{\mathcal{T}_{k,j}}(\hat{u}_k, u_{k,j}) = \sqrt{\mu} (\|\hat{u}_k - u_{k,j}\| + |u_{k,j} - \mathcal{I}_{\mathcal{T}_{k,j}} u_{k,j}|_{1,\Omega})$$

whence

$$\delta \geq \sqrt{\mu} \hat{C}_L \varepsilon_k.$$

Since $\|u - \hat{u}_k\| \leq \frac{\sqrt{\mu}}{2} \hat{C}_L \varepsilon_k$, we deduce that \hat{u}_k is an $\sqrt{\mu} \hat{C}_L \varepsilon_k$ -approximation of order s to u . Therefore, there exist an admissible mesh \mathcal{T}_δ such that

$$\mathbb{E}_{\mathcal{T}_\delta}(\hat{u}_k, u_{\mathcal{T}_\delta}^0) \leq \delta \quad \#\mathcal{T}_\delta \lesssim |u|_{\mathcal{A}_s}^{\frac{1}{s}} \delta^{-\frac{1}{s}}$$

where $u_{\mathcal{T}_\delta}^0 \in \mathbb{V}_{\mathcal{T}_\delta}^0$ because $\mathcal{A}_s = \mathcal{A}_s^0$. This implies

$$\|\hat{u}_k - u_{\mathcal{T}_\delta}^0\| = \mathbb{E}_{\mathcal{T}_\delta}(\hat{u}_k, u_{\mathcal{T}_\delta}) \leq \delta.$$

In order to compare with $\mathcal{T}_{k,j}$ we consider the overlay $\mathcal{T}_* = \mathcal{T}_{k,j} \oplus \mathcal{T}_\delta$, which satisfies

$$\#\mathcal{T}_* \leq \#\mathcal{T}_{k,j} + \#\mathcal{T}_\delta - \#\mathcal{T}_0.$$

Consider now $u_{\mathcal{T}_*}^0 \in \mathbb{V}_{\mathcal{T}_*}^0$, the Galerkin solution on the subspace of continuous piecewise linears $\mathbb{V}_{\mathcal{T}_*}^0$. Exploiting the monotonicity

$$\|\widehat{u}_k - u_{\mathcal{T}_*}^0\| \leq \|\widehat{u}_k - u_{\mathcal{T}_\delta}^0\|,$$

because \mathcal{T}_* is a refinement of \mathcal{T}_δ , we see that

$$\mathbb{E}_{\mathcal{T}_*}(\widehat{u}_k, u_{\mathcal{T}_*}^0) = \|\widehat{u}_k - u_{\mathcal{T}_*}^0\| \leq \|\widehat{u}_k - u_{\mathcal{T}_\delta}^0\| \leq \delta = \sqrt{\mu} \mathbb{E}_{\mathcal{T}_{k,j}}(\widehat{u}_k, u_{k,j}). \quad (8.38)$$

Applying Lemma 8.18 (optimal marking) to \mathcal{T}_* and $\mathcal{T}_{k,j}$ we infer that the refined set $R_{k,j} = R_{\mathcal{T}_{k,j} \rightarrow \mathcal{T}_*}$ satisfies Dörfler marking with parameter $0 < \theta < \theta^*$ and stabilization constant $\gamma > \gamma_*$. In addition,

$$\#R_{k,j} = \#\mathcal{T}_* - \#\mathcal{T}_{k,j}.$$

Since our Dörfler marking involves a minimal set $\mathcal{M}_{k,j}$, we deduce

$$\#\mathcal{M}_{k,j} \leq \#R_{k,j} \leq \#\mathcal{T}_\delta - \#\mathcal{T}_0 \lesssim |u|_{\mathcal{A}_s}^{\frac{1}{s}} \delta^{-\frac{1}{s}} \lesssim |u|_{\mathcal{A}_s}^{\frac{1}{s}} \varepsilon_k^{-\frac{1}{s}}.$$

This concludes the proof. \square

Corollary 8.20 (complexity of GALERKIN). *If $u \in \mathcal{A}_s$ and $\omega = \frac{\sqrt{\mu}\widehat{C}_L}{2D}$, the number of marked elements \mathcal{M}_k within a call to GALERKIN satisfies*

$$\#\mathcal{M}_k \leq JC_0 |u|_{\mathcal{A}_s}^{\frac{1}{s}} \varepsilon_k^{-\frac{1}{s}}.$$

Proof. Use that $\#\mathcal{M}_k = \sum_{j=0}^{J_k-1} \#\mathcal{M}_{k,j}$ and the previous lemma. \square

8.5 Quasi-optimality of AVEM

We finally address the quasi-optimality of the 2-loop method AVEM, by proving the announced bound (8.1).

Theorem 8.21 (quasi-optimality of AVEM). *Let Assumptions 8.5, 8.10, and 8.11 hold true. Then, there exist constants $\theta_*, \omega_* < 1$ and $\gamma_* \geq 1$ such that for all $\theta < \theta_*$, $\omega < \omega_*$, and $\gamma \geq \gamma_*$ there holds*

$$\|u - u_k\| \leq C(u, \mathcal{D}) (\#\mathcal{T}_k)^{-s} \quad 1 \leq k \leq K+1,$$

where $0 < s = \min\{s_u, s_{\mathcal{D}}\} = \min\{s_u, s_A, s_c, s_f\} \leq \frac{1}{2}$.

Proof. We know that the number of marked elements $N_k(u)$ within GALERKIN satisfies

$$N_k(u) \lesssim |u|_{\mathcal{A}_{s_u}}^{\frac{1}{s_u}} \varepsilon_k^{-\frac{1}{s_u}}$$

with $s_u \leq \frac{1}{2}$. Moreover, by Assumption 8.11 the number of marked elements $N_k(\mathcal{D})$ within DATA satisfies

$$N_k(\mathcal{D}) \lesssim |\mathcal{D}|_{\mathbb{A}_{s_{\mathcal{D}}}}^{\frac{1}{s_{\mathcal{D}}}} \varepsilon_k^{-\frac{1}{s_{\mathcal{D}}}}$$

with $s_{\mathcal{D}} \leq \frac{1}{2}$. Upon termination, **DATA** and **GALERKIN** give

$$\begin{aligned}\|u - \hat{u}_k\| &\leq D\omega\varepsilon_k = D\frac{\sqrt{\mu}\hat{C}_L}{2D}\varepsilon_k < \hat{C}_U\varepsilon_k, \\ \|\hat{u}_k - u_{k+1}\| &\leq \hat{C}_U\eta_{\mathcal{T}_{k+1}}(u_{k+1}, \mathcal{D}_k) \leq \hat{C}_U\varepsilon_k,\end{aligned}$$

because $\mu < 1$. This implies by triangle inequality

$$\|u - u_{k+1}\| \leq 2\hat{C}_U\varepsilon_k. \quad (8.39)$$

In addition, the total number of marked elements in the j -th loop of **AVEM** is

$$N_j(\mathcal{D}) + N_j(u) \leq C_1(|u|_{\mathcal{A}_{s_u}}^{\frac{1}{s_u}} + |\mathcal{D}|_{\mathbb{A}_{s_{\mathcal{D}}}}^{\frac{1}{s_{\mathcal{D}}}})\varepsilon_j^{-\frac{1}{s}}.$$

Therefore, the total amount of elements created by k loops of **AVEM**, besides those in \mathcal{T}_0 , obey the expression

$$\#\mathcal{T}_{k+1} - \#\mathcal{T}_0 \leq C_0 \sum_{j=0}^{k-1} (N_j(\mathcal{D}) + N_j(u)) \leq C_0 C_1 (|u|_{\mathcal{A}_{s_u}}^{\frac{1}{s_u}} + |\mathcal{D}|_{\mathbb{A}_{s_{\mathcal{D}}}}^{\frac{1}{s_{\mathcal{D}}}}) \sum_{j=0}^{k-1} \varepsilon_j^{-\frac{1}{s}}.$$

Since $\varepsilon_0 = 1$, $\varepsilon_j = 2^{-j}$ and

$$\sum_{j=0}^{k-1} (2^{-\frac{1}{s}})^j \leq \frac{1}{1 - 2^{-1/s}}$$

we deduce

$$\#\mathcal{T}_{k+1} - \#\mathcal{T}_0 \leq C(|u|_{\mathcal{A}_{s_u}}^{\frac{1}{s_u}} + |\mathcal{D}|_{\mathbb{A}_{s_{\mathcal{D}}}}^{\frac{1}{s_{\mathcal{D}}}})\varepsilon_k^{-\frac{1}{s}} \quad (8.40)$$

with $C = \frac{C_0 C_1}{1 - 2^{-1/s}}$. Since the first refined mesh satisfies $\#\mathcal{T}_1 \geq c_0 \#\mathcal{T}_0$ for some $c_0 > 1$, it holds $\#\mathcal{T}_{k+1} \leq \frac{c_0}{c_0 - 1}(\#\mathcal{T}_{k+1} - \#\mathcal{T}_0)$. Combining this with (8.40) and (8.39) yields the thesis. \square

Remark 8.22. The thresholds θ_*, ω_* play no role in Proposition 6.5 but are critical in Theorem 8.21. The former takes care of the gap between C_L and C_U in the a posteriori bounds (4.20), and is well documented in the optimality analysis of AFEMs [12, 19, 18, 20]. The latter guarantees that the perturbation error (6.15) is much smaller than ε_k and enables **GALERKIN** to learn the regularity of u from $\hat{u}_{\hat{\mathcal{T}}_k}$ [11, 20].

9 Data approximation: cardinality properties

In this section, we provide sufficient regularity conditions for data $\mathcal{D} = (A, c, f)$ to belong to the approximation classes introduced earlier and present algorithms for their approximation.

9.1 Greedy algorithm: definition and performance

We start with a constructive approximation estimate for a generic function $g : \Omega \rightarrow \mathbb{R}$ of class $W_p^s(\Omega)$ and next apply it to \mathcal{D} .

Let $1 \leq p, q \leq \infty$, $0 \leq s \leq 1$ be so that

$$\text{sob}(W_p^s(\Omega)) = s - \frac{2}{p} \geq \text{sob}(L^q(\Omega)) = 0 - \frac{2}{q},$$

whence

$$s - \frac{2}{p} + \frac{2}{q} \geq 0. \quad (9.1)$$

Let $E \in \mathcal{T}$ be a generic element, and let

$$g_E := \frac{1}{|E|} \int_E g$$

denote the mean value of g on E . Polynomial approximation theory yields

$$\|g - g_E\|_{L^q(E)} \lesssim h_E^{s-2/p+2/q} |g|_{W_p^s(E)}.$$

In view of the application to \mathcal{D} , it is convenient to consider the weighted $L^q(E)$ -norm instead, namely for $0 \leq t \leq 1$

$$\zeta_{\mathcal{T}}(E; g) := h_E^t \|g - g_E\|_{L^q(E)} \lesssim h_E^r |g|_{W_p^s(E)}, \quad \text{with } r := t + s - \frac{2}{p} + \frac{2}{q}. \quad (9.2)$$

Given a tolerance $\delta > 0$, we consider the algorithm

```
[ $\mathcal{T}$ ] = GREEDY( $\mathcal{T}, \delta$ )
while  $\mathcal{M} = \{E \in \mathcal{T} : \zeta_{\mathcal{T}}(E; g) > \delta\} \neq \emptyset$ 
   $\mathcal{T} = \text{REFINE}(\mathcal{T}, \mathcal{M})$ 
end while
return( $\mathcal{T}$ )
```

The following properties are valid for the global weighted error

$$\zeta_{\mathcal{T}}(g) := \left(\sum_{E \in \mathcal{T}} \zeta_{\mathcal{T}}^q(E; g) \right)^{\frac{1}{q}},$$

with the usual interpretation $\zeta_{\mathcal{T}}(g) := \max_{E \in \mathcal{T}} \zeta_{\mathcal{T}}(E; g)$ for $q = \infty$.

Proposition 9.1 (performance of GREEDY). *If $r > 0$, then GREEDY terminates in a finite number of steps. The output partition \mathcal{T} satisfies the estimates*

$$\zeta_{\mathcal{T}}(g) \leq \delta (\#\mathcal{T})^{\frac{1}{q}}, \quad (9.3)$$

$$\delta \lesssim |g|_{W_p^s(\Omega)} (\#\mathcal{T} - \#\mathcal{T}_0)^{-\frac{1}{q} - \frac{t+s}{2}}. \quad (9.4)$$

Remark 9.2 (error decay in GREEDY). Assuming $\#\mathcal{T} \geq c_0 \#\mathcal{T}_0$ for some $c_0 > 1$, and concatenating (9.3) and (9.4) yields

$$\zeta_{\mathcal{T}}(g) \lesssim |g|_{W_p^s(\Omega)} (\#\mathcal{T})^{-\frac{t+s}{2}}. \quad (9.5)$$

Proof of Proposition 9.1. We proceed in several steps.

- (i) *Termination.* Since $r > 0$, **GREEDY** stops in finite steps k , producing k subsequent refinements $\mathcal{T}_1, \dots, \mathcal{T}_k$ of \mathcal{T} . Upon termination, it holds $\zeta_{\mathcal{T}_k}(E; g) \leq \delta$ for all $E \in \mathcal{T}_k$, whence (9.3) follows.
- (ii) *Counting.* Let $\mathcal{M} = \mathcal{M}_0 \cup \dots \cup \mathcal{M}_{k-1}$ be the set of marked elements. We reorganize \mathcal{M} by size: let \mathcal{P}_j be the set of elements $E \in \mathcal{M}$ such that

$$2^{-(j+1)} \leq |E| < 2^{-j}, \quad \text{namely} \quad 2^{-\frac{j+1}{2}} \leq h_E < 2^{-\frac{j}{2}}.$$

Since **REFINE** uses bisection, the elements of \mathcal{P}_j are disjoint, whence

$$2^{-(j+1)} \# \mathcal{P}_j \leq |\Omega| \quad \text{i.e.,} \quad \# \mathcal{P}_j \leq |\Omega| 2^{j+1}.$$

On the other hand, $E \in \mathcal{P}_j$ (with $E \in \mathcal{T}_i$ for some i) implies

$$\delta < \zeta_{\mathcal{T}_i}(E; g) \lesssim h_E^r |g|_{W_p^s(E)} \leq 2^{-\frac{jr}{2}} |g|_{W_p^s(E)}.$$

In view of the summability of the right-hand side, we now accumulate these inequalities in the ℓ^p norm

$$\delta^p \# \mathcal{P}_j \lesssim 2^{-\frac{jrp}{2}} |g|_{W_p^s(\Omega)}^p.$$

This gives an alternative bound

$$\# \mathcal{P}_j \lesssim \delta^{-p} 2^{-\frac{jrp}{2}} |g|_{W_p^s(\Omega)}^p.$$

- (iii) *Summing up.* Adding over j we obtain

$$\# \mathcal{M} = \sum_j \# \mathcal{P}_j \lesssim \sum_{j \leq j_0} |\Omega| 2^{j+1} + \sum_{j > j_0} \delta^{-p} 2^{-\frac{jrp}{2}} |g|_{W_p^s(\Omega)}^p,$$

where j_0 corresponds to the crossover of the two series, namely

$$|\Omega| 2^{j_0+1} \simeq \delta^{-p} 2^{-\frac{j_0 rp}{2}} |g|_{W_p^s(\Omega)}^p.$$

This implies

$$2^{j_0(1+\frac{rp}{2})} \simeq |\Omega|^{-1} |g|_{W_p^s(\Omega)}^p \delta^{-p},$$

and

$$1 + \frac{rp}{2} = 1 + \frac{p}{2} \left(t + s - \frac{2}{p} + \frac{2}{q} \right) = \frac{p}{2}(t + s) + \frac{p}{q} = p w, \quad \text{with} \quad w := \frac{1}{2}(t + s) + \frac{1}{q}.$$

We thus deduce

$$2^{j_0} \simeq |\Omega|^{-\frac{1}{pw}} |g|_{W_p^s(\Omega)}^{\frac{1}{w}} \delta^{-\frac{1}{w}}$$

and the two series amount to the same sum

$$\# \mathcal{M} \lesssim |\Omega|^{1-\frac{1}{pw}} |g|_{W_p^s(\Omega)}^{\frac{1}{w}} \delta^{-\frac{1}{w}}.$$

(iv) *Complexity.* Apply finally the estimate (5.5) that controls the number of elements in \mathcal{T}_k in terms of \mathcal{M} :

$$\#\mathcal{T}_k - \#\mathcal{T}_0 \lesssim \#\mathcal{M} \lesssim |\Omega|^{1-\frac{1}{pw}} |g|_{W_p^s(\Omega)}^{\frac{1}{w}} \delta^{-\frac{1}{w}}.$$

This in turn yields

$$\delta \lesssim |\Omega|^{w-\frac{1}{p}} |g|_{W_p^s(\Omega)} (\#\mathcal{T}_k - \#\mathcal{T}_0)^{-w},$$

which is the asserted inequality (9.4) in view of the definition of w . This concludes the proof. \square

We now apply Proposition 9.1 to data $\mathcal{D} = (A, c, f)$, starting with A . In this case, we have

$$t = 0, \quad q = \infty, \quad r = s - \frac{2}{p} > 0.$$

This allows for $s = 1, p > 2$ (i.e., $A \in (W_p^1(\Omega))^{2 \times 2}$), or $s > 0, p = \infty$ (i.e., $A \in (W_\infty^s(\Omega))^{2 \times 2} = (C^{0,s}(\bar{\Omega}))^{2 \times 2}$, the space of Hölder-continuous tensor fields of exponent s).

Corollary 9.3 (approximation of A). *If $A \in (W_p^s(\Omega))^{2 \times 2}$ with $0 < s \leq 1$ and $p > \frac{2}{s}$, then*

$$\|A - \hat{A}_{\mathcal{T}}\|_{L^\infty(\Omega)} \lesssim |\Omega|^{\frac{s}{2}-\frac{1}{p}} |A|_{W_p^s(\Omega)} (\#\mathcal{T})^{-\frac{s}{2}}. \quad (9.6)$$

Thus, A belongs to the approximation class $\mathbb{A}_{\frac{s}{2}}$, and the GREEDY algorithm provides a quasi-optimal approximation of A .

We next consider the reaction term c , for which we have

$$t = 1, \quad q = \infty, \quad r = s - \frac{2}{p} + 1 > 0.$$

The latter inequality is surely satisfied if condition (9.1) holds. Thus, we may take $s = 1, p = 2$ (i.e., $c \in H^1(\Omega)$), or $0 \leq s \leq 1, p = \infty$ (i.e., $c \in W_\infty^s(\Omega)$).

Corollary 9.4 (approximation of c). *If $c \in W_p^s(\Omega)$ with $0 \leq s \leq 1$ and $p \geq \frac{2}{s}$, then*

$$\|h(c - \hat{c}_{\mathcal{T}})\|_{L^\infty(\Omega)} \lesssim |\Omega|^{\frac{1+s}{2}-\frac{1}{p}} |c|_{W_p^s(\Omega)} (\#\mathcal{T})^{-\frac{1+s}{2}}. \quad (9.7)$$

Thus, c belongs to the approximation class $\mathbb{C}_{\frac{1+s}{2}}$, and the GREEDY algorithm provides a quasi-optimal approximation of c .

We conclude with the forcing term f , for which we have

$$t = 1, \quad q = 2, \quad r = s - \frac{2}{p} + 2 > 0.$$

Again, the latter inequality is implied by (9.1). Admissible cases are $0 \leq s \leq 1, p = 2$ (i.e., $f \in H^s(\Omega)$), or $s = 1, p = 1$ (i.e., $f \in W_1^1(\Omega)$).

Corollary 9.5 (approximation of f). *If $f \in W_p^s(\Omega)$ with $0 \leq s \leq 1$ and $p \geq \frac{2}{s+1}$, then*

$$\|h(f - \hat{f}_{\mathcal{T}})\|_{L^2(\Omega)} \lesssim |\Omega|^{\frac{s}{2}+1-\frac{1}{p}} |f|_{W_p^s(\Omega)} (\#\mathcal{T})^{-\frac{1+s}{2}}. \quad (9.8)$$

Thus, f belongs to the approximation class $\mathbb{F}_{\frac{1+s}{2}}$, and the GREEDY algorithm provides a quasi-optimal approximation of f .

Remark 9.6 (rates of convergence). We see that the most critical data term is A , whose approximation error decays, according to (9.6), with rate $-\frac{s}{2}$ ($0 < s \leq 1$) provided $A \in W_p^s(\Omega)$. If $s = 1$, $p > 2$, we get the best possible rate $-\frac{1}{2}$.

On the other hand, data c and f lead to a rate $-\frac{1+s}{2} < -\frac{1}{2}$ for any regularity $c, f \in W_p^s(\Omega)$ with $0 < s \leq 1$. This is observed in the numerical experiments of Sect. 10. If instead, c and f have minimal regularity for our AVEM to make sense, namely $c \in L^\infty(\Omega)$, $f \in L^2(\Omega)$, then the convergence rates are $-\frac{1}{2}$ for both data (i.e., $s = 0$).

9.2 A pseudo-greedy strategy for f

Since the local error estimators $\zeta_{\mathcal{T}}(E; f) = h_E \|f - \hat{f}\|_{L^2(E)}$ accumulate in ℓ^2 , the threshold δ of GREEDY is not directly related to the desired tolerance ε . In fact, all $\zeta_{\mathcal{T}}(E; f)$ could be rather small relative to ε and yet $\zeta_{\mathcal{T}}(f) = \|\mathbf{h}(f - \hat{f})\|_{L^2(\Omega)} > \frac{1}{3}\varepsilon$. A practical choice is $\delta = \max_{T \in \mathcal{T}} \zeta_{\mathcal{T}}(E; f)$, but the ensuing algorithm is inefficient. We propose a minor modification of GREEDY with similar properties as Dörfler's algorithm that hinges on the maximum strategy. We describe the algorithm for a generic function $f \in W_p^s(\Omega)$ in the general setting presented at the beginning of this section, then we restrict the result to the forcing f of Corollary 9.5.

Given $\theta \in (0, 1)$ and a tolerance $\delta > 0$, consider the algorithm

```
[ $\mathcal{T}$ ] = P-GREEDY( $\mathcal{T}, \delta$ )
while  $\zeta_{\mathcal{T}}(f) > \delta$ 
   $\mathcal{M} = \{E \in \mathcal{T} : \zeta_{\mathcal{T}}(E; f) \geq \theta \max_{E' \in \mathcal{T}} \zeta_{\mathcal{T}}(E'; f)\}$ 
   $\mathcal{T} = \text{REFINE}(\mathcal{T}, \mathcal{M})$ 
end while
return( $\mathcal{T}$ )
```

The following statement is the counterpart of Proposition 9.1 and Remark 9.2 for P-GREEDY.

Proposition 9.7 (performance of P-GREEDY). *Let r be defined in (9.2), and suppose $r > 0$. Then P-GREEDY terminates in a finite number of steps. The output partition \mathcal{T} satisfies the estimates*

$$\zeta_{\mathcal{T}}(f) \leq \delta \quad \text{and} \quad \zeta_{\mathcal{T}}(f) \lesssim |f|_{W_p^s(\Omega)} (\#\mathcal{T})^{-\frac{t+s}{2}}. \quad (9.9)$$

Proof. Since the proof is similar to that of Proposition 9.1, we only report the new ingredients. Let $\mathcal{T}_1, \dots, \mathcal{T}_k$ be the sequence of refinements produced by P-GREEDY, and $\mathcal{M}_1, \dots, \mathcal{M}_k$ be the sequence of marked elements, with $\mathcal{M} = \mathcal{M}_1 \cup \dots \cup \mathcal{M}_k$. Set

$$\mu_i := \max\{\zeta_{\mathcal{T}_i}(E; f) : E \in \mathcal{T}_i\} \quad (1 \leq i \leq k) \quad \text{and} \quad \mu := \mu_{k-1}.$$

Then, it holds

$$\zeta_{\mathcal{T}_k}(f) \leq \delta < \zeta_{\mathcal{T}_{k-1}}(f) \leq \mu (\#\mathcal{T}_{k-1})^{\frac{1}{q}} \leq \mu (\#\mathcal{T}_k)^{\frac{1}{q}}. \quad (9.10)$$

On the other hand, since REFINE does not increase the element estimators, one has $\mu_i \geq \mu$ for any i , whence

$$\zeta_{\mathcal{T}_i}(E, f) \geq \theta \mu_i \geq \theta \mu \quad \forall E \in \mathcal{M}_i, \quad \forall i.$$

Let us introduce the partition of \mathcal{M} into disjoint subsets \mathcal{P}_j as in the proof of Proposition 9.1. If $E \in \mathcal{P}_j$, denoting by i the index such that $E \in \mathcal{M}_i$, we get

$$\theta \mu \leq \zeta_{\mathcal{T}_i}(E, f) \lesssim h_E^r |f|_{W_p^s(\Omega)} \leq 2^{-\frac{ir}{2}} |f|_{W_p^s(\Omega)},$$

whence

$$\#\mathcal{P}_j \lesssim \theta^{-p} \mu^{-p} 2^{-\frac{irp}{2}} |f|_{W_p^s(\Omega)}^p.$$

As in the proof of Proposition 9.1, this yields

$$\mu \lesssim |f|_{W_p^s(\Omega)} (\#\mathcal{T}_k - \#\mathcal{T}_0)^{-w},$$

and we conclude using (9.10) and the bound $\#\mathcal{T}_k \geq c_0 \#\mathcal{T}_0$ for $c_0 > 1$. \square

If the forcing $f \in W_p^s(\Omega)$ with $0 \leq s \leq 1$ and $p \geq \frac{2}{s+1}$, as in Corollary 9.5, then (9.9) reads $\zeta_{\mathcal{T}}(f) \lesssim |f|_{W_p^s(\Omega)} (\#\mathcal{T})^{-\frac{1+s}{2}}$, i.e, **P-GREEDY** provides a quasi-optimal approximation of f with convergence rate $-\frac{1+s}{2}$. In particular, if $f \in L^2(\Omega)$, then $\zeta_{\mathcal{T}}(f) \lesssim \|f\|_{L^2(\Omega)} (\#\mathcal{T})^{-\frac{1}{2}}$.

10 Numerical results

In this section we present a numerical experiment to confirm the convergence and optimality properties of the 2-step algorithm **AVEM**. We consider problem (2.1) in the L-shaped domain $\Omega = (-1, 1)^2 \setminus ([0, 1] \times [-1, 0])$, with diffusion term $A = aI$, where

$$a(x, y) = 1 + \exp(-50((x + 0.5)^2 + (y + 0.5)^2)) + \exp(-50((x + 0.5)^2 + (y - 0.5)^2)),$$

and reaction term

$$c(x, y) = 1 + \exp(-50((x + 0.5)^2 + y^2)) + \exp(-50(x^2 + (y - 0.5)^2));$$

note that the Gaussians in the definition of a and c have the same intensity but are located in different places within Ω (see Figures 2 and 3). The load term f and the Dirichlet boundary conditions are chosen in accordance with the analytical solution

$$u(x, y) = r^{\frac{2}{3}} \sin(2\alpha/3) + \exp(-1000((x - 0.5)^2 + (y - 0.5)^2)),$$

where (r, α) are the polar coordinates around the origin. Notice that the exact solution u is singular at the reentrant corner: it belongs to the Sobolev spaces $H(\Omega)^{\frac{5}{3}-\epsilon}$ with $\epsilon > 0$ and $W_p^2(\Omega)$ with $p > 1$. It also exhibits a rapid transition of order $10^{-3/2}$ around the point $(0.5, 0.5)$ due to the presence of a very narrow Gaussian. The three Gaussians are meant to test the performance of the module **DATA**.

We utilize the following parameters in the numerical test

$$\gamma = 1, \quad \Lambda = 10, \quad \theta_{\text{Dörfler}} = 0.5, \quad \omega = 1, \quad \theta_{\text{p-greedy}} = \text{sqrt}(0.75), \quad \text{tol} = 0.125,$$

where γ is parameter of the **dofi-dofi** stabilization (4.2), Λ is the bound for the global index of non-conforming partitions in Definition 3.2, $\theta_{\text{Dörfler}}$ is the Dörfler marking parameter (5.4), ω is the safety input parameter of **DATA**, $\theta_{\text{p-greedy}}$ is the pseudo-greedy marking parameter

(6.14), and `tol` is the target tolerance of **AVEM**. We implement algorithm **AVEM** with a uniform structured triangular mesh \mathcal{T}_0 with diameter $h = 0.125$ and initial tolerance $\epsilon_0 = 1$.

In order to estimate the VEM error between the exact solution u and the VEM solution $u_{\mathcal{T}}$, we consider the computable H^1 -like error quantity:

$$\text{H}^1\text{-error} := \frac{|u - \Pi_{\mathcal{T}}^{\nabla} u_{\mathcal{T}}|_{1,\mathcal{T}}}{|u|_{1,\Omega}}.$$

In Fig. 1 (left) we display the estimator $\eta_{\mathcal{T}}(u_{\mathcal{T}}, \mathcal{D})$, the data error $\zeta_{\hat{\mathcal{T}}}(\mathcal{D})$ and the $\text{H}^1\text{-error}$ obtained with algorithm **AVEM**. In Fig. 1 (right) we exhibit the data error $\zeta_{\hat{\mathcal{T}}}(\mathcal{D})$ and the addends $\zeta_{\hat{\mathcal{T}}}(A)$, $\zeta_{\hat{\mathcal{T}}}(c)$, $\zeta_{\hat{\mathcal{T}}}(f)$ (cf. (6.12) and (6.11)). Notice that the number of iterations of the algorithm **AVEM** is $K = \log_2(\epsilon_0/\text{tol}) = 3$.

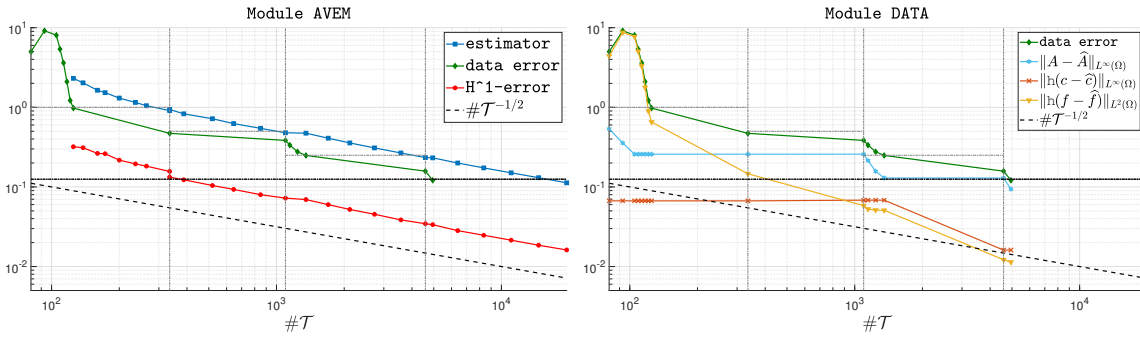


Figure 1: Left: estimator $\eta_{\mathcal{T}}(u_{\mathcal{T}}, \mathcal{D})$, data error $\zeta_{\hat{\mathcal{T}}}(\mathcal{D})$, $\text{H}^1\text{-error}$ obtained with the algorithm **AVEM**. Right: data error $\zeta_{\hat{\mathcal{T}}}(\mathcal{D})$, tensor error $\zeta_{\hat{\mathcal{T}}}(A)$, reaction error $\zeta_{\hat{\mathcal{T}}}(c)$, load error $\zeta_{\hat{\mathcal{T}}}(f)$, obtained with the algorithm **AVEM**. In both figures the optimal decay is indicated by the dashed line with slope -0.5 .

The predictions of Theorem 5.2 (contraction property of **GALERKIN**) are confirmed: both the estimator $\eta_{\mathcal{T}}(u_{\mathcal{T}}, \mathcal{D})$ and the $\text{H}^1\text{-error}$ converge to zero and the decay rate reaches asymptotically the theoretical optimal value $\#T^{-1/2}$; this corresponds to $s = 1/2$ in Theorem 8.21 (optimality of **AVEM**). Concerning data approximation, we observe from Fig. 1 (right) that $\zeta_{\hat{\mathcal{T}}}(\mathcal{D})$ decays with rate $\#T^{-1/2}$ dictated by $\zeta_{\hat{\mathcal{T}}}(A)$, as predicted by Corollary 9.3, while $\zeta_{\hat{\mathcal{T}}}(c)$ and $\zeta_{\hat{\mathcal{T}}}(f)$ exhibit a faster decay rate. This is due to regularity of (c, f) beyond $L^\infty(\Omega) \times L^2(\Omega)$, as predicted by Corollaries 9.4 and 9.5. We finally notice from Fig. 1 that the module **DATA** is active for all k except $k = 1$ because $\zeta_{\mathcal{T}_1}(\mathcal{D}) < \epsilon_1$.

In order to highlight the different level of approximation of data $\mathcal{D} = (A, c, f)$ required by **AVEM**, we display in Figs. 2, 3 and 4 the graphs of the piecewise constant approximations $\hat{\mathcal{D}} = (\hat{A}, \hat{c}, \hat{f})$ with respect to the mesh $\hat{\mathcal{T}}_K$ (left), and of the continuous piecewise linear counterparts with respect to the mesh \mathcal{T}_{K+1} (right). Since the Gaussians in a and c are located in non-overlapping subregions of Ω , it is possible to see that **AVEM** imposes a much finer resolution of a than of c in both meshes $\hat{\mathcal{T}}_K$ and \mathcal{T}_{K+1} ; this is due to the extra factor h in the definition (6.11) of $\zeta_{\hat{\mathcal{T}}}(c)$.

Finally in Figs. 6, 7 and 8 we compare the grids $\hat{\mathcal{T}}_K$ and \mathcal{T}_{K+1} generated by the modules **DATA** and **GALERKIN** upon termination of **AVEM**. The heat map on the rightmost pictures shows, for each element $E \in \mathcal{T}_{K+1}$, the number of newest-vertex bisections needed to create E starting

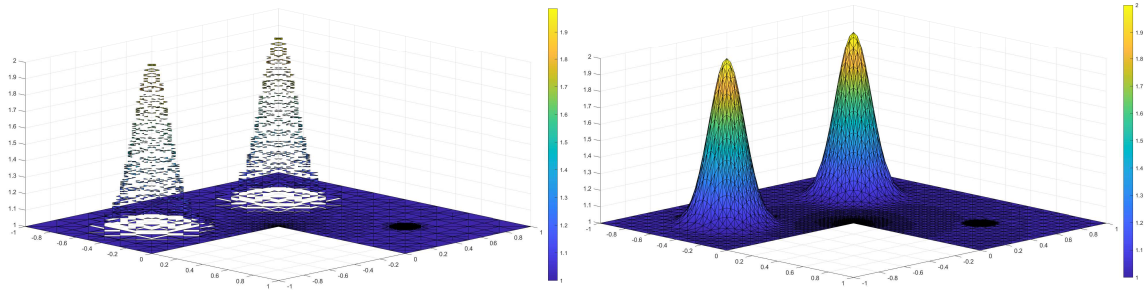


Figure 2: Left: graph of the piecewise constant approximation \hat{a} of a (w.r.t. $\hat{\mathcal{T}}_K$). Right: graph of the piecewise linear interpolant of a (w.r.t. \mathcal{T}_{K+1}).

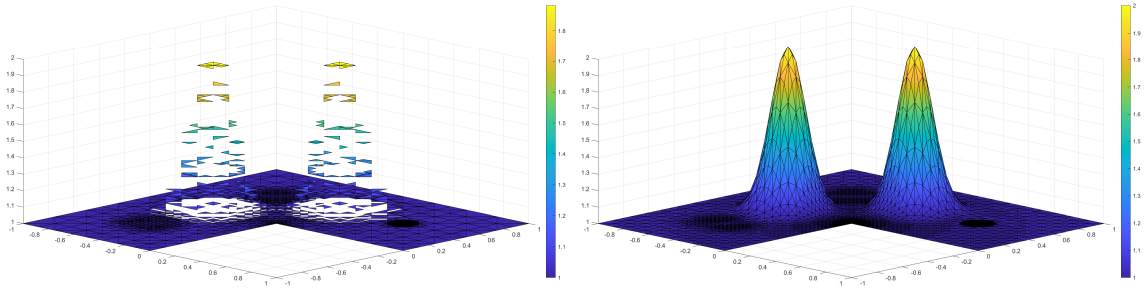


Figure 3: Left: graph of the piecewise constant approximation \hat{c} of c (w.r.t. $\hat{\mathcal{T}}_K$). Right: graph of the piecewise linear interpolant of c (w.r.t. \mathcal{T}_{K+1}). Notice much coarser resolution than in Fig. 2.

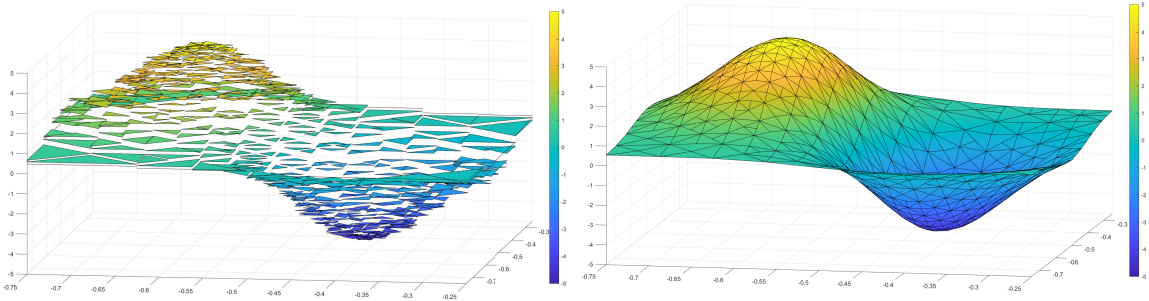


Figure 4: Zoom to $(-0.75, -0.25)^2$ for the load term f . Left: graph of the piecewise constant approximation \hat{f} of f (w.r.t. $\hat{\mathcal{T}}_K$). Right: graph of the piecewise linear interpolant of f (w.r.t. \mathcal{T}_{K+1}).

from $\hat{\mathcal{T}}_K$, according to the colorbar in Fig. 5. The number of nodes `N.vertices` and elements `N.elements` are

$$\begin{aligned} \text{N.vertices}(\hat{\mathcal{T}}_K) &= 5030, \\ \text{N.vertices}(\mathcal{T}_{K+1}) &= 19676, \end{aligned}$$

$$\begin{aligned} \text{N.elements}(\hat{\mathcal{T}}_K) &= 9236, \\ \text{N.elements}(\mathcal{T}_{K+1}) &= 37244. \end{aligned}$$



Figure 5: Colorbar for the heat map in Figures 6, 7 and 8.

Furthermore, the number of polygons in $\hat{\mathcal{T}}_K$ (elements with more than three vertices) is 730: 723 quadrilaterals, 2 pentagons, 5 hexagons; the number of polygons in \mathcal{T}_{K+1} is 1920: 1908 quadrilaterals, 16 pentagons, 4 hexagons. In Fig. 7 we plot a zoom to $(0.35, 0.65)^2$ of the meshes $\hat{\mathcal{T}}_K$ and \mathcal{T}_{K+1} . We highlight for both meshes the presence of hexagons in this subregion. Moreover, looking at the vertices having maximum global index λ sitting on the hexagons, we realize that the global indices are $\Lambda_{\hat{\mathcal{T}}_K} = 2$ and $\Lambda_{\mathcal{T}_{K+1}} = 3$. It is worth noting that the threshold $\Lambda = 10$ is never reached by **AVEM**; therefore, the condition of Λ -admissibility is not restrictive in practice. We further notice that the Gaussian in $(0.5, 0.5)$ associated with f is sufficiently resolved by **DATA**. In Fig. 8 we present a zoom to $(-10^{-2}, 10^{-2})^2$ to examine mesh refinement at the origin. We see that the mesh \mathcal{T}_{K+1} exhibits a rather strong grading at the reentrant corner, in accordance with the singularity of the exact solution. Elements in \mathcal{T}_{K+1} in this region need up to five newest-vertex bisection refinements relative to $\hat{\mathcal{T}}_K$.

We close this section with the following observation. From Figs. 6, 7 and 8 it can be appreciated how the presence of hanging nodes allows for quite abrupt and ‘steep’ refinements where needed in order to approximate the data and the solution singularity. In this respect, a direct comparison with AFEM in terms of generated meshes can be found in [8]. Such numerical results suggest that, although as shown in Remark 8.6 the approximation classes of **AVEM** and **AFEM** are the same, this added flexibility may be an important asset in adaptivity, especially in situations with more complex geometry. This aspect is worth further investigation, but is not within the scopes of the present contribution.

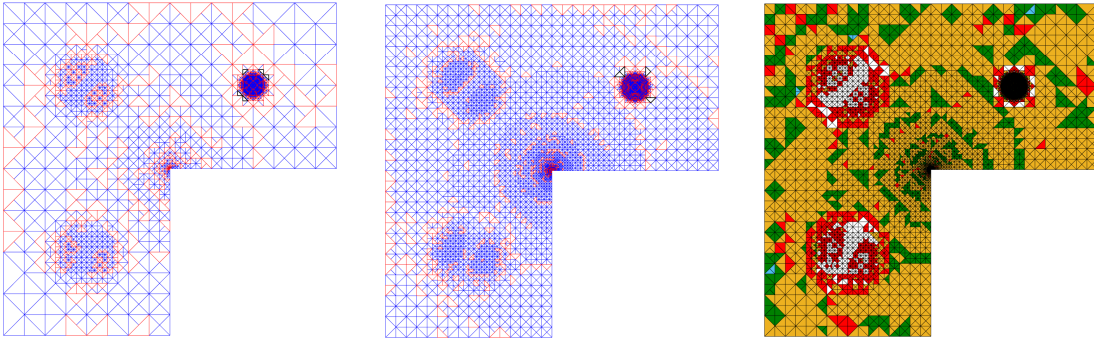


Figure 6: Left: final grid $\hat{\mathcal{T}}_K$ generated by **DATA**. Middle: final grid \mathcal{T}_{K+1} generated by **GALERKIN**. Mesh elements having more than three vertices (polygons) are drawn in red. Right: heat map representing for each $E \in \mathcal{T}_{K+1}$ the number of newest-vertex bisection needed to generate E starting from the mesh $\hat{\mathcal{T}}_K$ (colorbar in Fig. 5).

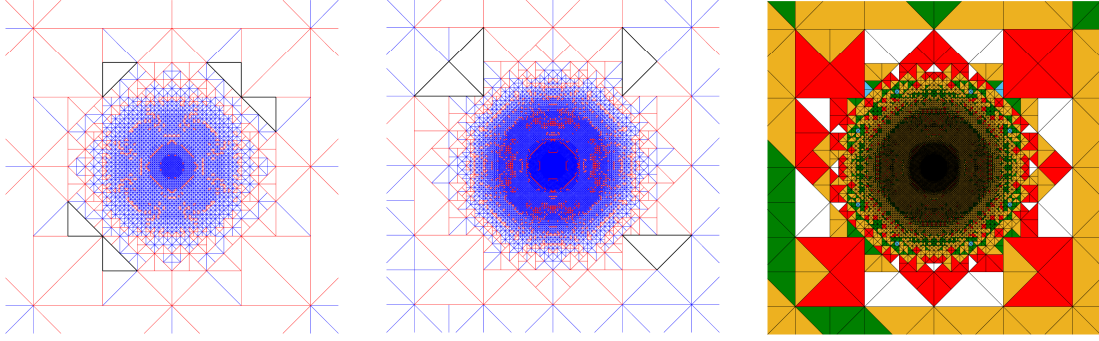


Figure 7: Zoom to $(0.35, 0.65)^2$ related to f . Left: final grid $\widehat{\mathcal{T}}_K$ generated by DATA. Middle: final grid \mathcal{T}_{K+1} generated by GALERKIN. Elements having more than three vertices (polygons) are drawn in red; elements drawn in black are hexagons. Right: heat map representing for each $E \in \mathcal{T}_{K+1}$ the number of newest-vertex bisection needed to generate E starting from the mesh $\widehat{\mathcal{T}}_K$ (colorbar in Fig. 5).

11 Λ -admissibility

Our theory of AVEM relies on the Λ -admissibility condition in Definition 3.2. Hereafter, we establish two results related to this concept: in Sect. 11.1, we show how to preserve Λ -admissibility while refining a mesh, whereas in Sect. 11.2 we prove that the overlay of two Λ -admissible meshes remains Λ -admissible.

11.1 Λ -admissible mesh refinement

In this section we introduce a constructive procedure that enforces Λ -admissibility at every stage of AVEM and study its complexity. If \mathcal{T} is a Λ -admissible refinement of \mathcal{T}_0 by newest-vertex bisection, the *level* of an element $E \in \mathcal{T}$, denoted by $\ell(E)$, is the number of successive bisections needed to generate E from \mathcal{T}_0 . Given $E \in \mathcal{T}$ marked for refinement, the procedure

$$[\mathcal{T}_*] = \text{CREATE_ADMISSIBLE_CHAIN}(\mathcal{T}, E, \Lambda)$$

generates a Λ -admissible refinement \mathcal{T}_* of \mathcal{T} upon bisecting E and at most $\ell(E) + 1$ other elements. To describe and analyze this procedure, we need some auxiliary notation and results.

Given any $E \in \mathcal{T}$, let us denote its newest vertex by $\mathbf{nv}(E)$, the edge opposite to $\mathbf{nv}(E)$ by $\mathbf{oe}(E)$, and the midpoint of $\mathbf{oe}(E)$ by $\mathbf{moe}(E)$. Furthermore, two elements $E', E'' \in \mathcal{T}$ are said *adjacent* if $e = E' \cap E''$ is an edge for at least one element, and are said *compatible* if they are adjacent and neither $\mathbf{nv}(E')$ nor $\mathbf{nv}(E'')$ belong to the line containing e (see Fig. 9, cases A and B).

Denote by \mathbb{T} the infinite tree obtained by successive bisections of the root partition \mathcal{T}_0 . The following result is well-known [10, 12, 18, 19, 21].

Lemma 11.1 (levels of elements sharing a full edge). *Assume that $E, E' \in \mathbb{T}$ share a full edge $e = E \cap E'$. Then*

$$|\ell(E) - \ell(E')| \leq 1.$$

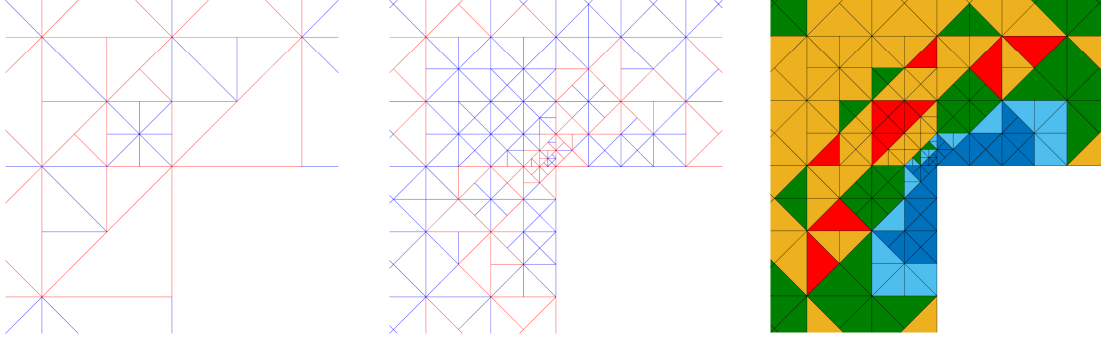


Figure 8: Zoom to $(-10^{-2}, 10^{-2})^2$ to examine the origin. Left: final grid $\widehat{\mathcal{T}}_K$ generated by DATA. Middle: final grid \mathcal{T}_{K+1} generated by GALERKIN. Elements having more than three vertices (polygons) are drawn in red. Right: heat map representing for each $E \in \mathcal{T}_{K+1}$ the number of newest-vertex bisection needed to generate E starting from the mesh $\widehat{\mathcal{T}}_K$ (colorbar in Fig. 5).

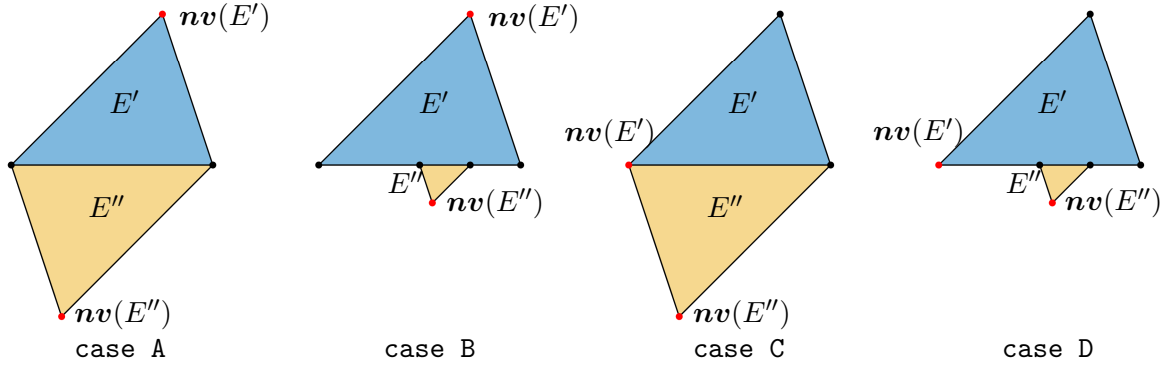


Figure 9: The elements E' and E'' are adjacent in cases A to D. They are compatible in cases A and B, and non-compatible in cases C and D.

Proof. If neither $nv(E)$ nor $nv(E')$ belong to e , or both $nv(E)$ and $nv(E')$ belong to e , then $\ell(E) = \ell(E')$. On the other hand, if $nv(E) \in e$ but $nv(E') \notin e$, then $\ell(E') = \ell(E) + 1$, since E' is generated by bisecting an element \widetilde{E} of the same level as E . \square

Lemma 11.2 (global index of a hanging node). *Consider an edge $[\mathbf{x}', \mathbf{x}']$ of the partition \mathcal{T} . If $\mathbf{x} \in \mathcal{H} \cap \text{int } e$ is generated by $m \geq 1$ bisections of e , then its global index $\lambda(\mathbf{x})$ satisfies*

$$\lambda(\mathbf{x}) = \max(\lambda(\mathbf{x}'), \lambda(\mathbf{x}'')) + m.$$

Proof. If $m = 1$, $\mathbf{x} = \mathbf{x}_M$ is the midpoint of e , and the formula is just the Definition 3.1 of global index. If $m > 1$, then \mathbf{x} is generated by bisecting some interval $[\mathbf{z}', \mathbf{z}'] \subset e$, and $\lambda(\mathbf{x}) = \max(\lambda(\mathbf{z}'), \lambda(\mathbf{z}'')) + 1$. Exactly one between $\mathbf{z}', \mathbf{z}''$ has been generated by $m - 1$ bisections, whereas the other one has been generated by less than $m - 1$ bisections. Hence, one concludes by induction. \square

Proposition 11.3 (reducing the global index of hanging nodes). *Let $\mathcal{H} \cap \text{int } e$ contain at least the midpoint \mathbf{x}_M of e . Assume that a bisection of some element in \mathcal{T} transforms \mathbf{x}_M into a proper node, and let λ_{new} denote the new global-index mapping of the nodes in $\mathcal{H} \cap \text{int } e$ after the bisection. Then there holds*

$$\lambda_{\text{new}}(\mathbf{x}) \leq \lambda(\mathbf{x}) - 1 \quad \forall \mathbf{x} \in \mathcal{H} \cap \text{int } e.$$

Proof. If $\mathbf{x} = \mathbf{x}_M$, then trivially $\lambda_{\text{new}}(\mathbf{x}) = 0 \leq \lambda(\mathbf{x}) - 1$. If $\mathbf{x} \in \mathcal{H} \cap \text{int } e$ is contained, say, in $(\mathbf{x}', \mathbf{x}_M)$ and has been generated by $m > 1$ successive bisections of e , then it is generated by $m - 1$ successive bisections of $[\mathbf{x}', \mathbf{x}_M]$. Thus, by Lemma 11.2

$$\begin{aligned} \lambda_{\text{new}}(\mathbf{x}) &\leq \max(\lambda_{\text{new}}(\mathbf{x}'), \lambda_{\text{new}}(\mathbf{x}_M)) + m - 1 \\ &= \max(\lambda(\mathbf{x}'), 0) + m - 1 = \lambda(\mathbf{x}') + m - 1 \\ &\leq \max((\lambda(\mathbf{x}'), \lambda(\mathbf{x}'')) + m - 1 = \lambda(\mathbf{x}) - 1. \end{aligned}$$

This gives the desired estimate. \square

The result just established is the motivation for the proposed refinement strategy. Indeed, it assures that in order to reduce the global index of a hanging node sitting on an edge, it is enough to transform the midpoint of the edge into a proper node.

The following remark will be useful in the sequel.

Remark 11.4 (facing element). Given a Λ -admissible mesh \mathcal{T} and $E \in \mathcal{T}$, let $\mathbf{x} = \mathbf{moe}(E)$ and suppose that $\lambda(\mathbf{x}) > \Lambda$. Then \mathbf{x} is not a node of \mathcal{T} , whence the edge $\mathbf{oe}(E)$ cannot contain any hanging node in its interior. We conclude that there exists a unique adjacent element $\tilde{E} \in \mathcal{T}$, $\tilde{E} \neq E$, such that $E \cap \tilde{E} = \mathbf{oe}(E)$. This element will be called the element *facing* E .

Given an element $E \in \mathcal{T}$ which has been marked for refinement, we are ready to identify those elements in \mathcal{T} that need be bisected with E in order to create a Λ -admissible refinement of \mathcal{T} .

Definition 11.5 (chain of elements to be refined). *Define by recurrence the chain of elements*

$$\mathcal{C}(E) = \{E_0, E_1, \dots, E_K\}$$

for some $K \geq 0$, as follows: set first $E_0 = E$ and, assuming to have defined E_k for $k \geq 0$, then

- (i) *if $\lambda(\mathbf{moe}(E_k)) \leq \Lambda$, set $K = k$ and stop;*
- (ii) *if $\lambda(\mathbf{moe}(E_k)) = \Lambda + 1$ and the facing element \tilde{E}_k is compatible with E_k , set $E_{k+1} = \tilde{E}_k$, $K = k + 1$ and stop;*
- (iii) *if $\lambda(\mathbf{moe}(E_k)) = \Lambda + 1$ and the facing element \tilde{E}_k is not compatible with E_k , set $E_{k+1} = \tilde{E}_k$ and continue.*

Lemma 11.6 (properties of the chain of refinement). *The chain $\mathcal{C}(E)$ has at most $K \leq \ell(E) + 1$ elements. Furthermore, the sequence of element levels $\{\ell(E_k)\}_{k=0}^K$ is not increasing.*

Proof. We claim that step (iii) in Definition 11.5 reduces the level by at least one. In fact, E_k coincides with or is a refinement of a triangle $E \in \mathbb{T}$ sharing with E_{k+1} a full edge; thus $\ell(E_k) \geq \ell(E)$. Such triangle E satisfies $\ell(E) = \ell(E_{k+1}) + 1$ according to Lemma 11.1, whence

$$\ell(E_{k+1}) = \ell(E) - 1 \leq \ell(E_k) - 1. \quad (11.1)$$

Therefore, for as long as case (iii) is active, i.e. for all $j < K$, we have $\ell(E_j) \leq \ell(E_0) - j$ and

$$0 \leq \ell(E_{K-1}) \leq \ell(E_0) - (K - 1),$$

which gives the first part of the Lemma. The monotonicity of $\{\ell(E_k)\}_{k=0}^K$ follows from (11.1) and the fact that $\ell(E_{K-1}) = \ell(E_K)$ in case (ii). \square

We are now ready to define the procedure

$$[\mathcal{T}_*] = \text{CREATE_ADMISSIBLE_CHAIN}(\mathcal{T}, E, \Lambda)$$

The partition \mathcal{T}_* is obtained from \mathcal{T} by refining only the elements in $\mathcal{C}(E)$. More precisely, starting from E_K , one goes traverses the chain backwards and, for $K \geq k \geq 1$, considers the cases

- if E_k and E_{k-1} are compatible, then E_k is bisected once (see Fig. 10, cases A or B);
- if E_k and E_{k-1} are not compatible, then E_k is bisected twice and, after the first bisection, the sibling that is facing E_{k-1} is further bisected (see Fig. 10, cases C or D);
- finally, $E_0 = E$ is bisected once.

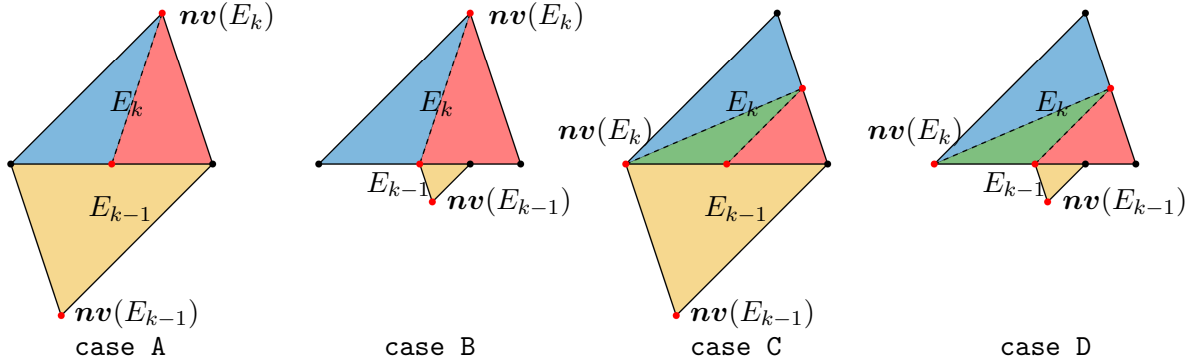


Figure 10: Two elements E_{k-1} and E_k in the chain $\mathcal{C}(E)$: E_{k-1} can be bisected in a Λ -admissible way, only after E_k is refined once (cases A and B), or twice (cases C and D)

Proposition 11.7 (properties of `CREATE_ADMISSIBLE_CHAIN`). *If \mathcal{T} is Λ -admissible, then the call $[\mathcal{T}_*] = \text{CREATE_ADMISSIBLE_CHAIN}(\mathcal{T}, E, \Lambda)$ bisects once or twice the elements of the chain $\mathcal{C}(E)$, whose cardinality is at most $\ell(E) + 1$, and produces a Λ -admissible mesh \mathcal{T}_* with E bisected once. Moreover, every element $E' \in \mathcal{T}_*$ generated by this call satisfies*

$$\ell(E') \leq \ell(E) + 1. \quad (11.2)$$

Proof. Let $\mathcal{C}(E) = \{E_k\}_{k=0}^K$ and observe that, for $k \geq 1$, one or two bisections of E_k convert the midpoint of the edge e of E_k shared with E_{k-1} into a proper node. Therefore, Proposition 11.3 (reducing the global index of hanging nodes) implies that the global indices of all interior nodes to e decrease by at least 1, and makes the bisection of E_{k-1} Λ -admissible as desired. This procedure creates \mathcal{T}_* upon partitioning at most $\ell(E) + 1$ elements, namely those of $\mathcal{C}(E)$, according to Lemma 11.6 (properties of the chain refinement).

Moreover, to prove (11.2) we take $k \geq 1$ and consider the following two mutually exclusive cases. If E_k and E_{k-1} are compatible, then E_k is replaced by two elements $E' \in \mathcal{T}_*$ of level

$$\ell(E') = \ell(E_k) + 1 \leq \ell(E) + 1,$$

according to Lemma 11.6. On the other hand, if E_k and E_{k-1} are not compatible, then E_k is replaced by one element of level $\ell(E_k) + 1$ and two elements $E' \in \mathcal{T}_*$ of level

$$\ell(E') = \ell(E_k) + 2 \leq \ell(E_{k-1}) + 1 \leq \ell(E) + 1$$

because of (11.1). Finally, the element $E_0 = E$ is replaced by two elements of level $\ell(E) + 1$. \square

In view of Proposition 11.7 a bound of the form $\#\mathcal{T}_* - \#\mathcal{T} \leq C_0$ with a universal constant C_0 is false because C_0 may depend on $\ell(E)$ in general. This obstruction to optimal complexity of **REFINE** was tackled by Binev, Dahmen and DeVore in their seminal paper [10], and further studied in [12, 18, 19, 21]. In fact, the cumulative effect of bisection on conforming meshes obeys the weaker, but yet optimal, equation (5.5). The extension of this to Λ -admissible non-conforming partitions is precisely guaranteed by the stated Theorem 5.1, whose proof follows.

*Proof of Theorem 5.1 (complexity of **REFINE**).* We follow [19, Section 6.3], which explains the basic ingredients to derive (5.5). It turns out that two crucial properties of **CREATE_ADMISSIBLE_CHAIN** as required. The first is (11.2). The second one relates the level of elements and their distance to E , namely

$$\text{dist}(E, E') \leq C 2^{\frac{\ell(E')}{2}} \quad \forall E' \in \mathcal{T}_* \setminus \mathcal{T};$$

such property is valid for bisection grids regardless of Λ -admissibility [19, Lemma 18]. This completes the proof. \square

11.2 Mesh Overlay and Λ -admissibility

Given two partitions \mathcal{T}_A and \mathcal{T}_B , denote by $\mathcal{T}_A \oplus \mathcal{T}_B$ the *overlay* of \mathcal{T}_A and \mathcal{T}_B , i.e., the partition whose associated tree is the union of the trees of \mathcal{T}_A and \mathcal{T}_B . The following property holds.

Proposition 11.8. *If \mathcal{T}_A and \mathcal{T}_B are Λ -admissible, then $\mathcal{T}_A \oplus \mathcal{T}_B$ remains Λ -admissible.*

Proof. Denote here by \mathcal{N} the set of all nodes obtained by newest-vertex bisection from the root partition \mathcal{T}_0 . Let $\mathcal{N}_0, \mathcal{N}_A, \mathcal{N}_B, \mathcal{N}_{A+B}$, resp., be the set of nodes of the partitions $\mathcal{T}_0, \mathcal{T}_A, \mathcal{T}_B, \mathcal{T}_A \oplus \mathcal{T}_B$, resp.. It is easily seen that for each $\mathbf{x} \in \mathcal{N} \setminus \mathcal{N}_0$ there exists a unique $\mathcal{B}(\mathbf{x}) = \{\mathbf{x}', \mathbf{x}''\} \subset \mathcal{N}$ such that \mathbf{x} is generated by the bisection of the segment $[\mathbf{x}', \mathbf{x}'']$. Furthermore, if $\mathbf{x} \in \mathcal{N}_{A+B}$ is a proper node of \mathcal{T}_A (of \mathcal{T}_B , resp.), then it is also a proper node of $\mathcal{T}_A \oplus \mathcal{T}_B$.

Let us denote by λ_A , λ_B , λ_{A+B} , resp., the global-index mappings defined on \mathcal{N}_A , \mathcal{N}_B , \mathcal{N}_{A+B} , resp.. It is convenient to extend the definition of λ_A and λ_B to the whole \mathcal{N}_{A+B} by setting

$$\lambda_A(\mathbf{x}) = +\infty \quad \text{if } \mathbf{x} \in \mathcal{N}_{A+B} \setminus \mathcal{N}_A, \quad \lambda_B(\mathbf{x}) = +\infty \quad \text{if } \mathbf{x} \in \mathcal{N}_{A+B} \setminus \mathcal{N}_B.$$

With these notations at hand, we are going to prove the inequality

$$\lambda_{A+B}(\mathbf{x}) \leq \min(\lambda_A(\mathbf{x}), \lambda_B(\mathbf{x})) \quad \forall \mathbf{x} \in \mathcal{N}_{A+B}, \quad (11.3)$$

from which the thesis immediately follows.

We proceed by induction on $k = \lambda_{A+B}(\mathbf{x})$, $\mathbf{x} \in \mathcal{N}_{A+B}$. If $k = 0$, the inequality is trivial since $\lambda_A(\mathbf{x}), \lambda_B(\mathbf{x}) \geq 0$. So suppose (11.3) hold up to some $k \geq 0$. If $\mathbf{x} \in \mathcal{N}_{A+B}$ satisfies $\lambda_{A+B}(\mathbf{x}) = k + 1 > 0$, then it is a hanging node of $\mathcal{T}_A \oplus \mathcal{T}_B$ by definition of global index, hence, it is a hanging node of \mathcal{T}_A or \mathcal{T}_B ; wlog, suppose it is a hanging node of \mathcal{T}_A . If \mathbf{x} is generated by the bisection of the segment $[\mathbf{x}', \mathbf{x}'']$, then again by definition of global index it holds

$$k + 1 = \lambda_{A+B}(\mathbf{x}) = \max(\lambda_{A+B}(\mathbf{x}'), \lambda_{A+B}(\mathbf{x}'')) + 1,$$

which implies

$$\lambda_{A+B}(\mathbf{x}') \leq k, \quad \lambda_{A+B}(\mathbf{x}'') \leq k.$$

By induction,

$$\lambda_{A+B}(\mathbf{x}') \leq \min(\lambda_A(\mathbf{x}'), \lambda_B(\mathbf{x}')), \quad \lambda_{A+B}(\mathbf{x}'') \leq \min(\lambda_A(\mathbf{x}''), \lambda_B(\mathbf{x}')),$$

from which we obtain

$$\lambda_{A+B}(\mathbf{x}) \leq \max(\lambda_A(\mathbf{x}'), \lambda_A(\mathbf{x}'')) + 1 = \lambda_A(\mathbf{x})$$

since \mathbf{x} is a hanging node of \mathcal{T}_A . On the other hand, either $\mathbf{x} \in \mathcal{N}_B$ or $\mathbf{x} \notin \mathcal{N}_B$. In the latter case, $\lambda_B(\mathbf{x}) = +\infty$, and (11.3) is proven. In the former case, necessarily \mathbf{x} is a hanging node of \mathcal{T}_B , hence as above

$$\lambda_{A+B}(\mathbf{x}) \leq \max(\lambda_B(\mathbf{x}'), \lambda_B(\mathbf{x}'')) + 1 = \lambda_B(\mathbf{x}),$$

and the thesis is proven. \square

12 Conclusions

This paper introduces and studies a two-step adaptive virtual element method (**AVEM**) of lowest order over triangular meshes with hanging nodes in 2d, which are treated as polygons. **AVEM** applies to linear symmetric elliptic problems with variable data. The main achievements of the paper can be summarized as follows:

- **AVEM** concatenates two modules, **DATA** and **GALERKIN**. The former approximates data by piecewise constants to a desired accuracy, while the latter handles the adaptive approximation of the problem with piecewise constant data, as described in [8]. **AVEM** converges (Proposition 6.5);

- *Complexity of GALERKIN*: the number of sub-iterations inside the call to **GALERKIN** at iteration k of **AVEM** is bounded independently of k (Proposition 7.3);
- *Complexity of DATA*: the module **DATA** is quasi-optimal in terms of accuracy versus mesh cardinality, under suitable regularity conditions on the data (Sect. 9);
- *Complexity of AVEM*: **AVEM** is quasi-optimal in terms of error decay versus degrees of freedom, for solutions and data belonging to appropriate approximation classes (Theorem 8.21);
- *Numerical experiments*: they illustrate the interplay between the modules **DATA** and **GALERKIN** and provide computational evidence of the optimality of **AVEM** (Section 10).
- *Mesh admissibility*: Section 11 designs a procedure to keep the global index of meshes uniformly bounded for all steps k , and proves its optimality in terms of degrees of freedom.

Although in Remark 8.6 we observed that, in the presence of a bound on the maximal index of hanging nodes, the equivalence classes of **AVEM** and **AFEM** are the same, the numerical results in Section 10 and in [8] suggest that the flexibility of VEM may lead to more efficient meshes in complex situations, at least in terms of the involved constants. A deeper investigation of this aspect at the theoretical level may require a more advanced VEM approach, for instance taking inspiration from the a-priori analysis in [4].

Acknowledgements

LBdV, CC and MV were partially supported by the Italian MIUR through the PRIN grants n. 201744KLJL and n. 20204LN5N5 (LBdV, MV) and n. 201752HKH8 (CC). RHN has been supported in part by NSF grant DMS-1908267. These supports are gratefully acknowledged. LBdV, CC, MV and GV are members of the INdAM research group GNCS.

References

- [1] B. Ahmad, A. Alsaedi, F. Brezzi, L. D. Marini, and A. Russo. Equivalent projectors for virtual element methods. *Comput. Math. Appl.*, 66(3):376–391, 2013.
- [2] L. Beirão da Veiga, C. Lovadina, and A. Russo. Stability analysis for the virtual element method. *Math. Mod.and Meth. in Appl. Sci.*, 27(13):2557–2594, 2017.
- [3] L. Beirão da Veiga and G. Manzini. Residual *a posteriori* error estimation for the virtual element method for elliptic problems. *ESAIM Math. Model. Numer. Anal.*, 49(2):577–599, 2015.
- [4] L. Beirão da Veiga and G. Vacca. Sharper error estimates for virtual elements and a bubble-enriched version. *SIAM J. Numer. Anal.*, 60(4):1853–1878, 2022.
- [5] L. Beirão da Veiga, F. Brezzi, A. Cangiani, G. Manzini, L. D. Marini, and A. Russo. Basic principles of virtual element methods. *Math. Models Methods Appl. Sci.*, 23(1):199–214, 2013.

- [6] L. Beirão da Veiga, F. Brezzi, L. D. Marini, and A. Russo. The Hitchhiker’s Guide to the Virtual Element Method. *Math. Models Methods Appl. Sci.*, 24(8):1541–1573, 2014.
- [7] L. Beirão da Veiga, F. Brezzi, L. D. Marini, and A. Russo. Virtual Element Method for general second-order elliptic problems on polygonal meshes. *Math. Models Methods Appl. Sci.*, 24(4):729–750, 2016.
- [8] L. Beirão da Veiga, C. Canuto, R. H. Nochetto, G. Vacca, and M. Verani. Adaptive vem: Stabilization-free a posteriori error analysis and contraction property. *in press on SINUM*, 2022.
- [9] S. Berrone and A. Borio. A residual *a posteriori* error estimate for the Virtual Element Method. *Math. Models Methods Appl. Sci.*, 27(8):1423–1458, 2017.
- [10] P. Binev, W. Dahmen, and R. DeVore. Adaptive finite element methods with convergence rates. *Numer. Math.*, 97(2):219–268, 2004.
- [11] A. Bonito, R. A. DeVore, and R. H. Nochetto. Adaptive finite element methods for elliptic problems with discontinuous coefficients. *SIAM J. Numer. Anal.*, 51(6):3106–3134, 2013.
- [12] A. Bonito and R. H. Nochetto. Quasi-optimal convergence rate of an adaptive discontinuous Galerkin method. *SIAM J. Numer. Anal.*, 48(2):734–771, 2010.
- [13] S. C. Brenner and L.Y. Sung. Virtual element methods on meshes with small edges or faces. *Math. Models Methods Appl. Sci.*, 28(7):1291–1336, 2018.
- [14] A. Cangiani, E. H. Georgoulis, T. Pryer, and O. J. Sutton. A posteriori error estimates for the virtual element method. *Numer. Math.*, 137(4):857–893, 2017.
- [15] C. Carstensen, M. Feischl, M. Page, and D. Praetorius. Axioms of adaptivity. *Comput. Math. Appl.*, 67(6):1195–1253, 2014.
- [16] A. Cohen, R. DeVore, and R. H. Nochetto. Convergence rates of AFEM with H^{-1} data. *Found. Comput. Math.*, 12(5):671–718, 2012.
- [17] W. Dörfler. A convergent adaptive algorithm for Poisson’s equation. *SIAM J. Numer. Anal.*, 33(3):1106–1124, 1996.
- [18] R. H. Nochetto, K. G. Siebert, and A. Veiser. Theory of adaptive finite element methods: an introduction. In *Multiscale, nonlinear and adaptive approximation*, pages 409–542. Springer, Berlin, 2009.
- [19] R. H. Nochetto and A. Veiser. Primer of adaptive finite element methods. In *Multiscale and adaptivity: modeling, numerics and applications*, volume 2040 of *Lecture Notes in Math.*, pages 125–225. Springer, Heidelberg, 2012.
- [20] R. Stevenson. Optimality of a standard adaptive finite element method. *Found. Comput. Math.*, 7(2):245–269, 2007.
- [21] R. Stevenson. The completion of locally refined simplicial partitions created by bisection. *Math. Comp.*, 77(261):227–241, 2008.